

# Learning PySpark: A Guide to Removing Spaces from DataFrame Column Names

Authored by  
**Mohammed loot**

November 11, 2025

## RECOMMENDED CITATION

Mohammed loot (2025). *Learning PySpark: A Guide to Removing Spaces from DataFrame Column Names*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=16652>

Working with large-scale data processing requires rigorous attention to detail, especially when managing the structure of a [DataFrame](#). One common challenge faced by data engineers using [PySpark](#) is dealing with inconsistent or poorly formatted column names, such as those containing spaces. While spaces are syntactically valid in many database systems, they often complicate querying, analysis, and subsequent data manipulation when using libraries like [PySpark](#) or integration tools that prefer camelCase or snake\_case nomenclature. This article provides a comprehensive guide on how to efficiently and programmatically remove or replace spaces within column names in a [PySpark DataFrame](#) using functional programming techniques.

## The Core Technique: Using PySpark's select and alias Methods

The most robust and idiomatic way to modify column names in [PySpark](#) involves iterating through the existing column list and applying a transformation using the **select** method in conjunction with the **alias** function. This approach ensures that the entire [DataFrame](#) structure is preserved while only the metadata (the column names) is updated. We leverage list comprehension for conciseness and efficiency, creating a new list of columns where each element is renamed dynamically.

To implement this transformation, we must import the necessary [functions](#) module, typically aliased as `F`. The method involves selecting each column (`F.col(x)`) and then applying a standard Python string manipulation technique (`x.replace(' ', '_')`) before assigning the new name using the [alias](#) method. This technique is highly scalable and recommended for production environments where consistency in naming conventions is paramount.

Below is the precise syntax required to iterate over the column names and replace all spaces with an underscore character (`_`). This transformation is typically executed early in the data pipeline to ensure cleaner access later:

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

```
# Define the transformation: replace all spaces in column names with underscores  
df_new = df.select()
```

Understanding this block is critical: `df.columns` provides a list of current column names. The list comprehension iterates through this list (`for x in df.columns`). For each column name `x`, we first apply the string method `x.replace(' ', '_')` to create the desired new name. Finally, `F.col(x).alias(...)` maps the original column data (retrieved by `F.col(x)`) to the new, cleaned column name using the [alias](#) operation. This powerful combination allows for extremely flexible and efficient renaming across large datasets.

## Practical Implementation: Setting Up the PySpark Environment and Sample Data

To demonstrate this functionality in a real-world scenario, we first need to establish a working [PySpark](#) environment and create a sample [DataFrame](#) that intentionally includes spaces in its column headers. We begin by initializing the [SparkSession](#), which is the required entry point for all functionality in Spark. This setup process ensures that we have the necessary context to perform data definition and manipulation tasks using the distributed computing framework.

Our example focuses on basketball statistics, where column names like "team name," "points scored," and "total assists" are descriptive but contain problematic spaces. These spaces would typically hinder direct SQL-like querying within Spark unless special quoting mechanisms are used, which is generally discouraged for simplicity and maintainability in production code.

The following code block outlines the standard procedure for defining the data, specifying the column names (including the spaces), creating the [DataFrame](#), and displaying the initial structure:

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

# Define sample data rows
data = ,
,
,
,
,
,
]

# Define column names containing spaces
columns =

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

# View the initial DataFrame structure
df.show()

+-----+-----+-----+
|team name|points scored|total assists|
+-----+-----+-----+
| Mavs| 18| 4|
```

```
| Nets| 33| 9|
| Hawks| 12| 7|
| Thunder| 15| 3|
| Lakers| 19| 2|
| Cavs| 24| 5|
| Magic| 28| 7|
+-----+-----+-----+
```

As clearly demonstrated in the output above, the column headers--"team name," "points scored," and "total assists"--all contain standard whitespace characters. Our primary objective in the following steps is to use the programmatic approach detailed earlier to clean these names, making them compliant with standard data warehousing and analytical practices such as `snake_case`.

### Example 1: Replacing Spaces with Underscores (Snake Case)

Replacing spaces with underscores is the most common and highly recommended practice when standardizing column names in data engineering. This convention, commonly known as **snake\_case**, improves readability significantly compared to completely concatenating words, while ensuring the column name remains a single token accessible without special delimiters. This is particularly beneficial when integrating with external systems or databases that natively utilize SQL syntax, which often struggles with embedded spaces.

We will now apply the core list comprehension technique using the [functions](#) module to generate a new [DataFrame](#), `df_new`, where the spaces are substituted by underscores. It is important to remember that this process is non-destructive; it creates a new DataFrame object with the desired schema rather than modifying the original in place, adhering to the immutability principles of Spark.

The code below implements the renaming logic and immediately displays the result, confirming the successful transformation:

```
from pyspark.sql import functions as F

# Apply the renaming operation: replace spaces with underscores
df_new = df.select()

# View the newly created DataFrame
df_new.show()

+-----+-----+-----+
|team_name|points_scored|total_assists|
+-----+-----+-----+
```

```
| Mavs| 18| 4|
| Nets| 33| 9|
| Hawks| 12| 7|
| Thunder| 15| 3|
| Lakers| 19| 2|
| Cavs| 24| 5|
| Magic| 28| 7|
+-----+-----+-----+
```

Upon reviewing the output, it is evident that the transformation was successful: "team name" is now `team_name`, "points scored" is `points_scored`, and "total assists" is `total_assists`. This standardized naming convention dramatically simplifies future operations, such as joining DataFrames or writing complex transformation logic in [PySpark](#), ensuring compliance with preferred coding standards.

## Example 2: Complete Removal of Spaces (Concatenation)

While replacing spaces with underscores is generally preferred for readability, there are scenarios--particularly when enforcing strict **camelCase** or adhering to specific non-delimited naming conventions--where the complete removal of spaces is desired. This results in concatenated, single-word column names, which can be less visually intuitive but equally valid from a programmatic standpoint. Although less common in large-scale data warehousing than `snake_case`, understanding this alternative method provides complete flexibility in column hygiene.

Achieving this outcome requires only a minor modification to the string replacement logic used previously. Instead of replacing the space character (' ') with an underscore ('\_'), we replace it with an empty string (''). This leverages the intrinsic power of Python's built-in string functions applied efficiently within the Spark context to achieve the desired effect--the space is simply removed without introducing any substitute character.

The syntax for complete space removal is shown below, followed by the resulting DataFrame display. Notice the change in the second argument of the `replace` function, which determines the substitution character:

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

```
# Remove all spaces in column names by replacing them with an empty string
df_new = df.select()
```

```
# View the new DataFrame structure
df_new.show()
```

```
+-----+-----+-----+
|teamname|pointsscored|totalassists|
+-----+-----+-----+
| Mavs| 18| 4|
| Nets| 33| 9|
| Hawks| 12| 7|
| Thunder| 15| 3|
| Lakers| 19| 2|
| Cavs| 24| 5|
| Magic| 28| 7|
+-----+-----+-----+
```

As anticipated, the column names are now entirely concatenated: `teamname`, `pointsscored`, and `totalassists`. While functional, it is important for developers to weigh the programmatic benefits against the potential decrease in code clarity when choosing this approach over using underscores or other delimiters, especially when dealing with complex datasets.

## Best Practices and Advanced Considerations

The techniques demonstrated utilize the standard Python string `replace` function within a [PySpark](#) context to manage column metadata. This specific function is highly efficient when dealing with simple, single-character replacements like spaces. However, production data cleaning often requires more complex transformations than just removing spaces. For instance, normalizing special characters, removing leading/trailing whitespace, converting to lower case, or handling punctuation simultaneously are common requirements.

For more advanced renaming tasks, such as removing multiple types of special characters (e.g., parentheses, hyphens, or multiple consecutive spaces), we recommend using Python's **Regular Expressions (Regex)** library, `re`, within the same list comprehension structure. A Regex pattern allows for much greater control over which characters are targeted for replacement or removal, ensuring complete standardization across varied input sources. Regardless of the complexity of the string operation, the core structure involving `df.select()` remains the preferred and most scalable pattern for bulk column renaming in Spark.

When designing data pipelines, developers should establish a strict naming convention early on. The **snake\_case** standard (`lower_case_with_underscores`) is generally the safest bet for compatibility across Spark, SQL databases, and general Python analysis tools. By applying this column cleaning process as one of the very first steps in any Extract, Transform, Load (ETL) pipeline, you preemptively eliminate a significant source of errors, simplify future joins, and improve the maintainability of your data models.

## Additional Resources

Mastering column manipulation is fundamental to effective data engineering in PySpark. The techniques shown here--leveraging list comprehension combined with the **select** and [alias](#) methods--provide the foundation for highly efficient schema management. To deepen your understanding of these and other related functions, consult the official documentation provided by Apache Spark for the PySpark SQL module. Exploring these resources will help you perform other common data transformation tasks efficiently.

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