

Learning PySpark: Joining DataFrames with Mismatched Column Names

Authored by
Mohammed looti

November 11, 2025

RECOMMENDED CITATION

Mohammed looti (2025). *Learning PySpark: Joining DataFrames with Mismatched Column Names*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=16647>

The process of integrating disparate datasets is fundamental to modern data analysis and engineering. When working with [PySpark](#), joining two or more [DataFrames](#) is a routine operation. However, a common challenge arises when the corresponding linking columns in the source DataFrames possess different names. Standard join syntax requires identical column names, which necessitates a preparatory transformation.

Fortunately, PySpark provides an elegant and efficient solution utilizing the [withColumn](#) function, which allows for the dynamic renaming of columns within the join operation itself. This technique ensures that the join condition is satisfied without permanently altering the structure of the source DataFrames prior to execution.

Understanding the Challenge of Data Joins in PySpark

In most relational database systems, including the architecture that underpins [Apache Spark](#), the **join operation** relies on matching values found in specified key columns. If DataFrame A contains a column named `user_ID` and DataFrame B contains the corresponding key under the name `customer_id`, a direct join attempt using the column names as keys will fail, as the system perceives them as two distinct attributes.

The challenge is magnified in large-scale distributed computing environments like Spark, where data lineage and transformation efficiency are critical. While one could manually rename the column in one DataFrame, execute the join, and then potentially drop or rename the resulting columns, this approach often leads to verbose and less readable code. Furthermore, it requires multiple sequential operations, potentially breaking the desired flow of transformations.

Our goal is to execute the join while preserving the original datasets' integrity and achieving maximum code conciseness. The key is to create a temporary, unified column name that both DataFrames can reference simultaneously during the join execution phase.

The Solution: Using the `withColumn` Transformation

The standard syntax to join two DataFrames (`df1` and `df2`) based on different column names (e.g., `team_id` in `df1` and `team_name` in `df2`) involves chaining the `withColumn` transformation directly into the join logic. This method is highly recommended for its efficiency and readability in PySpark data pipelines.

The `withColumn` function is primarily used to add a new column or replace an existing column based on an expression. We strategically use it here to map the disparate key columns to a single, temporary column name--conventionally named `id` or `join_key`--within the scope of the join operation.

The syntax below illustrates this powerful approach. Notice how the operations are executed inline, ensuring that the temporary column exists only for the duration of the join calculation:

```
df3 = df1.withColumn('id', col('team_id')).join(df2.withColumn('id', col('team_name')), on='id')
```

This concise line of code performs several critical steps simultaneously, enabling a clean and effective join across mismatched keys.

Step-by-Step Implementation of the Join Logic

To fully understand the mechanism, let us break down the sequence of operations performed by the snippet above. This is crucial for debugging and optimizing complex transformations in PySpark.

The operation begins with the first DataFrame, `df1`, where we use `withColumn` to duplicate the value from the original key column (`team_id`) into a newly created column named `id`. This new column now serves as the standardized joining key for `df1`.

First, it renames the `team_id` column from `df1` to `id` by creating a new column with the same value. Then, it performs the same renaming/duplication process on `df2`, mapping the `team_name` column to the new `id` column.

Lastly, it executes the **join operation** between the two transformed DataFrames, specifying the standardized `id` column as the common key (`on='id'`).

This approach ensures that the resulting DataFrame, `df3`, contains the columns from both original DataFrames, including the redundant original keys (`team_id` and `team_name`) and the new `id` column used for matching. The following practical example demonstrates how this syntax is applied in a real-world scenario involving sports team data.

Detailed Practical Example Setup

To illustrate this technique, we will create two distinct [DataFrames](#). Imagine we have two separate data sources tracking basketball statistics. The first source uses a unique team identifier (`team_ID`), while the second uses the full team name (`team_name`). Despite the different naming conventions, the values within the columns are identical, allowing for a successful linkage.

We begin by initializing the Spark Session and defining our input data structures. Note the distinct column names between `df1` and `df2`, which sets up our joining challenge.

Suppose we have the following DataFrame named `df1`, containing points data linked by `team_ID`:

```

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

#define data
data = ,
,
,
,
,
,
]

#define column names
columns =

#create dataframe using data and column names
df1 = spark.createDataFrame(data, columns)

#view dataframe
df1.show()

```

```

+-----+-----+
|team_ID|points|
+-----+-----+
| Mavs| 18|
| Nets| 33|
| Lakers| 12|
| Kings| 15|
| Hawks| 19|
| Wizards| 24|
| Magic| 28|
+-----+-----+

```

Next, we define our second DataFrame, **df2**, which contains assist data, but uses the column name `team_name` for the identifier. This difference in naming is the core issue we must resolve using the inline column renaming technique.

And suppose we have another DataFrame named **df2**:

```

#define data
data = ,

```

```
,
,
,
,
,
]

#define column names
columns =

#create dataframe using data and column names
df2 = spark.createDataFrame(data, columns)

#view dataframe
df2.show()
```

```
+-----+-----+
|team_name|assists|
+-----+-----+
| Hawks| 4|
| Wizards| 5|
| Raptors| 5|
| Kings| 12|
| Mavs| 7|
| Nets| 11|
| Magic| 3|
+-----+-----+
```

Executing the Join and Analyzing the Result

With our two DataFrames defined, we can now apply the previously discussed syntax to perform an inner join. This join will link rows where the value in `df1.team_ID` matches the value in `df2.team_name`, using the temporary `id` column as the bridge.

We use the following syntax to perform the **inner join** between these two DataFrames by renaming the team columns from each DataFrame to `id` and then joining on values from the `id` column. It is important to remember that because we are using an inner join (the default join type in PySpark), only records present in both DataFrames will be included in the output. Notice, for instance, that 'Lakers' (only in `df1`) and 'Raptors' (only in `df2`) will be excluded.

```
#join df1 and df2 on different column names
```

```
df3 = df1.withColumn('id', col('team_id')).join(df2.withColumn('id', col('team_name')), on='id')
```

```
#view resulting DataFrame
```

```
df3.show()
```

```
+-----+-----+-----+-----+-----+
| id|team_ID|points|team_name|assists|
+-----+-----+-----+-----+-----+
| Hawks| Hawks| 19| Hawks| 4|
| Kings| Kings| 15| Kings| 12|
| Magic| Magic| 28| Magic| 3|
| Mavs| Mavs| 18| Mavs| 7|
| Nets| Nets| 33| Nets| 11|
| Wizards|Wizards| 24| Wizards| 5|
+-----+-----+-----+-----+-----+
```

The result, `df3`, successfully merges the points and assists data based on the matching team names. Critically, because the join was executed using the temporary `id` column, the output DataFrame now contains five columns: the new `id` column, the original `team_ID` and `points` columns from `df1`, and the original `team_name` and `assists` columns from `df2`. We have successfully joined the two DataFrames into one DataFrame based on matching values in the new `id` column.

While this result is correct, the presence of three columns representing the exact same entity (`id`, `team_ID`, and `team_name`) is redundant and can complicate subsequent analysis or storage. This leads us to the final step of structuring the output.

Streamlining Output: Using `select()` After the Join

In most data processing pipelines, the end goal is a clean, normalized dataset that includes only necessary columns. After performing a join using temporary columns, it is best practice to use the `select` function to filter the resulting DataFrame down to the required attributes, effectively dropping the unnecessary intermediate columns.

The `select` function allows us to specify exactly which columns from the joined DataFrame should be retained. We can choose to keep the standardized `id` column and the desired metrics (`points` and `assists`), while discarding the redundant original key columns (`team_ID` and `team_name`).

The following syntax demonstrates how to chain the `select` transformation immediately after the join operation. This creates a much cleaner output DataFrame, optimized for consumption by

reporting tools or subsequent ETL stages:

#join df1 and df2 on different column names

```
df3 = df1.withColumn('id', col('team_id')).join(df2.withColumn('id', col('team_name')), on='id')
.select('id', 'points', 'assists')
```

#view resulting DataFrame

```
df3.show()
```

```
+-----+-----+-----+
| id|points|assists|
+-----+-----+-----+
| Hawks| 19| 4|
| Kings| 15| 12|
| Magic| 28| 3|
| Mavs| 18| 7|
| Nets| 33| 11|
| Wizards| 24| 5|
+-----+-----+-----+
```

Notice that only the **id**, **points**, and **assists** columns are shown in the final joined DataFrame. This streamlined output is generally the preferred format for the final stages of a data transformation pipeline. This robust, single-line chained command is an excellent example of idiomatic PySpark programming, minimizing intermediate variables and maximizing efficiency in distributed data processing.

Further Resources and Advanced Concepts

Mastering joins on mismatched columns is a critical skill for any PySpark developer. This technique is often preferable to using the more complex SQL expression syntax within the join clause when dealing with simple key mismatches.

For those looking to expand their knowledge of PySpark transformations, it is highly recommended to explore additional join types (such as left, right, and full outer joins), which can be specified easily using the `how` parameter in the `join` function. Understanding how to manage schema evolution and optimize column operations, particularly in scenarios involving millions or billions of records, is the next step in becoming an expert in distributed data processing.

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

How to perform various types of joins in PySpark (Left, Right, Full).

Using window functions for complex aggregations.
Optimizing data partitioning for improved join performance.