

# Learning PySpark: Performing Left Joins with Multiple Columns

Authored by  
**Mohammed loot**

November 10, 2025

## RECOMMENDED CITATION

Mohammed loot (2025). *Learning PySpark: Performing Left Joins with Multiple Columns*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=16483>

## Understanding Joins in Distributed Data Processing

In the modern landscape of big data and [distributed computing](#), efficiently combining massive datasets is a core responsibility of any data engineer. Frameworks like [PySpark](#)--the Python API for Apache Spark--are specifically designed to handle these integration challenges at scale. When data is partitioned across multiple nodes, establishing accurate relationships between records often requires more than just a simple primary key match. This necessitates the use of a multi-column join, or [Composite Key](#) join, which is fundamental for ensuring data integrity and successful data enrichment across the entire cluster. Mastering the technique of joining [DataFrames](#) based on multiple criteria is crucial for anyone operating in this environment.

The complexity arises because the join operation must be executed reliably across all partitions simultaneously. Unlike traditional SQL databases where the data resides centrally, PySpark must coordinate the shuffle and comparison of data across worker nodes. Therefore, the definition of the join condition must be explicit and unambiguous. A multi-column join allows engineers to define a precise linkage, perhaps matching on a combination of 'customer ID' and 'transaction date,' ensuring that only truly corresponding records are merged. This precision minimizes the risk of producing Cartesian products or inaccurate data pairings, which can be catastrophic in large-scale analysis.

This article focuses specifically on the [Left Join](#), which is perhaps the most common join type used in data preparation pipelines. The purpose of the [Left Join](#) is protective: it guarantees that every record from the primary (left) [DataFrame](#) is preserved in the resulting output. Even if no corresponding match is found in the secondary (right) DataFrame, the left record remains, with the columns from the right side being filled with **null** values. This behavior is essential for auditing and tracing unmatched records, which is a vital component of any robust data quality process.

### The Mechanics of the Left Join and Composite Keys

When working with complex business data, defining a unique entity often requires combining several attributes. This combination forms what is known as a [Composite Key](#). For instance, in an organizational dataset, matching records might require comparing not just an employee ID, but also their department code and the fiscal year. In PySpark, executing a join based on a [Composite Key](#) demands a specific syntax where the conditions are treated as a collective set of Boolean expressions that must all evaluate to **true** for a match to occur.

The standard SQL approach to a [Left Join](#) is translated into the PySpark **.join()** method. The mechanism ensures that for every row in the left [DataFrame](#), Spark iterates through the partitions of the right DataFrame, checking the composite condition. If a match is found, the rows are combined. If no match is found, the data from the left side is carried forward, and the fields corresponding to the right side are populated with the special marker **null**.

Understanding the distinction between joining on a list of column names (which Spark handles implicitly if names are identical) versus joining on explicit Boolean expressions is vital for multi-key joins. When column names differ between the DataFrames--as is often the case when integrating data from disparate sources--we must use explicit column-to-column comparisons. This explicit definition, using syntax like `df1 == df2`, gives the data engineer precise control over the linkage and is mandatory when dealing with heterogeneous schemas.

## Implementing the Core PySpark Syntax for Multi-Key Joins

The power of [PySpark](#) lies in its expressive syntax, which translates complex data integration logic into concise code. To execute a multi-column **left join**, we leverage the `.join()` method and define the composite matching logic within the `on` parameter. Critically, this parameter accepts a list of individual equality conditions, where each condition specifies a pairing between a column in the left DataFrame and a column in the right DataFrame.

The required structure involves calling `df1.join(df2, on=, how='left')`. The list of conditions within the `on` argument acts as a logical AND operation: the record linkage only succeeds if **all** equality checks within the list are satisfied simultaneously. This rigorous requirement is precisely what defines the [Composite Key](#) match. Furthermore, specifying `how='left'` explicitly instructs the [distributed computing](#) engine to prioritize the retention of all records from the first specified DataFrame (`df1`).

The canonical syntax for performing a left join using two linking columns, `col1` and `col2`, where the column names might differ between the two source [DataFrames](#), is demonstrated below. This pattern is robust and handles situations where schema naming conventions are not synchronized:

```
df_joined = df1.join(df2, on=, how='left')
```

This code snippet effectively creates a new DataFrame, `df_joined`, which preserves the entirety of `df1`. Records from `df2` are only appended if and only if the values in `df2.col1` and `df2.col2` perfectly align with the values found in `df1.col1` and `df1.col2`, respectively. This explicit definition is key to managing data complexity in a scalable manner.

## Preparing Sample DataFrames for Demonstration

To effectively demonstrate the mechanics of this multi-key [Left Join](#), we must first establish a functional [SparkSession](#) and create two sample [DataFrames](#) that reflect real-world data integration scenarios. We will purposefully design the schemas to have differing column names for the join keys, forcing the use of the explicit Boolean expression syntax, and include mismatched records to observe the crucial behavior of the **left join**.

Our first DataFrame, **df1**, represents our primary dataset (the left side). It contains core statistics, such as player points, uniquely identified by a combination of **team** and **pos** (position). This DataFrame must retain all its records in the final output. Note the inclusion of the record 'B', 'G', 14, which will intentionally have no match in the secondary DataFrame, providing a clear test case for the **null** filling behavior.

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

```
#define data
```

```
data1 = ,
```

```
,
```

```
,
```

```
]
```

```
#define column names
```

```
columns1 =
```

```
#create dataframe using data and column names
```

```
df1 = spark.createDataFrame(data1, columns1)
```

```
#view dataframe
```

```
df1.show()
```

```
+----+----+-----+
|team|pos|points|
+----+----+-----+
| A| G| 18|
| A| F| 22|
| B| F| 19|
| B| G| 14|
+----+----+-----+
```

Our second DataFrame, **df2**, acts as the secondary dataset, providing supplementary information (like assists). Crucially, its columns are named **team\_name** and **position**, distinct from **df1**. Furthermore, **df2** contains data for Team C, which will be ignored entirely by the **left join**, demonstrating that only records matching the left side are considered for inclusion. This setup ensures a robust demonstration of explicit multi-key joining and the retention property of the [Left Join](#).

```
#define data
```

```

data2 = ,
,
,
,
]

#define column names
columns2 =

#create dataframe using data and column names
df2 = spark.createDataFrame(data2, columns2)

#view dataframe
df2.show()

+-----+-----+-----+
|team_name|position|assists|
+-----+-----+-----+
| A| G| 4|
| A| F| 9|
| B| F| 8|
| C| G| 6|
| C| F| 5|
+-----+-----+-----+

```

## Executing and Analyzing the Multi-Column Left Join

With our source [DataFrames](#) prepared, the next step is to execute the powerful multi-key [Left Join](#). We must define the composite condition carefully to link records only when both the team and the position identifiers align perfectly. Failure to match on even one component of the [Composite Key](#) will result in a non-match, preserving the left record while introducing **null** values for the right-side columns.

The join condition list explicitly defines the linkage: we equate **df1.team** with **df2.team\_name**, and simultaneously equate **df1.pos** with **df2.position**. By passing this list to the **on** parameter and specifying **how='left'**, we initiate the distributed computation across the cluster, ensuring the correct data relationship is maintained based on our composite criteria.

```

#perform left join
df_joined = df1.join(df2, on=, how='left')

```

```
#view resulting DataFrame
df_joined.show()

+----+----+-----+-----+-----+
|team|pos|points|team_name|position|assists|
+----+----+-----+-----+-----+
| A| G| 18| A| G| 4|
| A| F| 22| A| F| 9|
| B| F| 19| B| F| 8|
| B| G| 14| null| null| null|
+----+----+-----+-----+-----+
```

Analyzing the output confirms the expected behavior of the **left join**. The first three rows successfully found corresponding records in **df2**, resulting in the **assists** data being correctly appended. However, the last record, Team B, Position G, which was present in **df1** but lacked a match in **df2**, remains fully preserved. Its corresponding columns from the right side (**team\_name**, **position**, and **assists**) are accurately filled with **null** values, confirming the integrity of the operation and the successful implementation of the composite join criteria.

## Post-Join Cleanup: Eliminating Redundant Columns

A typical consequence of using explicit Boolean expressions for joining, especially when column names differ (e.g., **df1.key1 == df2.key2**), is that the resulting joined [DataFrame](#) contains duplicate key columns. In our example, we have both **team** and **team\_name**, and **pos** and **position**. While these columns hold identical values where a match occurred, retaining all four is inefficient, adding unnecessary overhead to storage and subsequent processing steps in a large-scale [distributed computing](#) environment.

The best practice following a **left join** is to streamline the output by dropping the redundant key columns originating from the right DataFrame (**df2**). We use the **.drop()** transformation, which can be chained directly onto the join operation or executed separately on the resulting DataFrame. This ensures that the final dataset is concise and optimized, relying solely on the keys from the primary (left) structure.

We apply the **.drop()** method to remove the right-side key columns, **team\_name** and **position**, producing the final, clean result set ready for analysis or further transformations:

```
#drop 'team_name' and 'position' columns from joined DataFrame
df_joined.drop('team_name', 'position').show()
```

```
+----+----+-----+-----+-----+
```

```
|team|pos|points|assists|
+----+----+-----+-----+
| A| G| 18| 4|
| A| F| 22| 9|
| B| F| 19| 8|
| B| G| 14| null|
+----+----+-----+-----+
```

This final output successfully integrates the data based on the composite criteria while maintaining a clean schema. Notice that the **null** value for the unmatched record (Team B, Position G) is correctly preserved in the **assists** column, demonstrating that the data integrity required by the **left join** was fully upheld throughout the process.

## Advanced Performance Tips and Best Practices

While correctly implementing the multi-column join syntax is essential, achieving high performance in a large-scale [PySpark](#) environment requires strategic optimization. Joins are notoriously resource-intensive operations in Spark because they often trigger a data shuffle, moving large amounts of data across the network between worker nodes. Minimizing this shuffle is key to efficiency.

One of the most effective optimization techniques is [Broadcast Joining](#). If the right DataFrame (**df2**) is small--typically less than 10MB to 100MB, depending on the cluster configuration--it can be explicitly broadcasted to all worker nodes. When Spark performs a broadcast join, the smaller DataFrame is copied entirely to the memory of every executor, allowing the join calculation to occur locally without requiring the costly network shuffle of the large left DataFrame. This strategy can yield dramatic performance improvements, often reducing join times from minutes to seconds.

Furthermore, consistency in data types is a non-negotiable best practice. Before executing any join operation, especially those involving [Composite Keys](#), developers must ensure that the corresponding join key columns in both [DataFrames](#) share identical data types (e.g., string to string, integer to integer). Mismatched types will cause the equality conditions to fail universally, resulting in an output where the right side is entirely populated by **null** values, leading to silent data errors that can be difficult to diagnose. Utilizing functions like **.withColumn()** and **.cast()** for type validation and coercion should be standard operating procedure in any robust data pipeline leveraging [SparkSession](#) functionality.

Mastering multi-column joins is essential for complex data integration tasks in [PySpark](#). By explicitly defining the composite key using a list of Boolean expressions within the **on** parameter, developers can accurately link records across distributed DataFrames, ensuring data integrity is

maintained at scale.

## **Additional Resources**

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