

# How to Replace Multiple Values in a PySpark Column

Authored by  
**stats writer**

January 2, 2026

## RECOMMENDED CITATION

stats writer (2026). *How to Replace Multiple Values in a PySpark Column*.

PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=110535>

Data cleaning and transformation are fundamental tasks in any data processing pipeline, and when working with big data environments, tools like [PySpark](#) are indispensable. A common requirement is the need to efficiently replace multiple specific values within a single column of a [DataFrame](#). While simple replacements might use basic mapping, complex or large-scale transformations require specialized, performance-optimized functions native to the Spark ecosystem. Understanding the various methods available for these transformations is crucial for maintaining code clarity and computational efficiency.

The core challenge in replacing multiple values simultaneously lies in implementing conditional logic across the entire column efficiently. Unlike standard Python libraries, PySpark operates on distributed data, meaning solutions must leverage vectorized operations rather than row-by-row processing. The initial text mentions a dedicated `replace()` function, which is useful for direct dictionary mapping. However, for nuanced conditional logic, particularly when dealing with strings or when retaining unmatched values is necessary, the combination of the `when()` and `otherwise()` functions provided by the `pyspark.sql.functions` module proves to be the most flexible and robust approach. This method allows users to define a series of specific conditions, analogous to `CASE WHEN` statements in SQL, ensuring powerful control over the data modification process.

In this guide, we will focus primarily on using the powerful `when()` and `otherwise()` chain to achieve precise, multiple value replacements within a PySpark DataFrame column. We will also briefly discuss the direct `replace()` method as an alternative. Mastering this technique is a cornerstone of effective data manipulation within the Spark environment, enabling complex data standardizations, categorizations, and corrections efficiently across massive datasets. The ability to handle these replacements programmatically ensures reproducibility and scalability in data engineering workflows.

## Method 1: Utilizing the `when()` and `otherwise()` Functions

The PySpark function `when()` is analogous to a multi-branch `if-then-else` structure. It allows data engineers to specify a condition and the resulting value if that condition is met. When chaining multiple `when()` clauses together, Spark evaluates them sequentially. The first condition that evaluates to true dictates the replacement value for that row. This chaining mechanism is perfect for implementing the replacement of multiple distinct values with corresponding new values in a single pass over the column data.

Crucially, the `when()` chain must be terminated by the `otherwise()` function. The `otherwise()` clause defines the default action for any row that does not meet any of the preceding conditions. In the context of value replacement, the most common practice is to use `otherwise(df.column_name)`, which dictates that if a value is not explicitly targeted for

replacement, it should retain its original value. This prevents unintended data loss or null values being introduced for categories that should remain untouched.

To implement this transformation, we use the `withColumn()` method of the `DataFrame`. The `withColumn()` operation is non-destructive; it returns a new `DataFrame` resulting from the modification. We specify the column to be modified (e.g., 'team') and then define the transformation logic using the chained `when()` and `otherwise()` functions. This approach ensures that we are leveraging Spark's optimized execution engine, which handles the conditional logic efficiently across all partitions of the data. This structured approach is highly readable and ensures that complex mapping logic is executed reliably.

## Detailed Syntax Breakdown for Conditional Replacement

The core syntax for implementing conditional replacements using `when()` requires importing the necessary functions from the PySpark library. Specifically, both the `when` function and the `col` function (or direct column referencing, as shown in the example below) are often needed. Below is the standard structure used to replace multiple values in a specific column of a PySpark `DataFrame`:

```
from pyspark.sql.functions import when
```

```
#replace multiple values in 'team' column
df_new = df.withColumn('team', when(df.team=='A', 'Atlanta')
    .when(df.team=='B', 'Boston')
    .when(df.team=='C', 'Chicago'))
    .otherwise(df.team))
```

This block of code demonstrates the concise power of the `when()` chain. Each `when()` clause takes two arguments: the condition (e.g., `df.team=='A'`) and the replacement value if the condition is true (e.g., `'Atlanta'`). The series of conditional checks is executed for every row. If a row's 'team' value is 'A', it immediately gets replaced with 'Atlanta', and the subsequent `when()` clauses are skipped for that specific row. This short-circuit evaluation is key to both performance and logical accuracy.

The final `otherwise(df.team)` ensures that any existing values in the 'team' column that do not match 'A', 'B', or 'C' are preserved without modification. For instance, if a row contained the value 'D', the previous `when()` conditions would all fail, and the `otherwise()` clause would be triggered, setting the new 'team' value equal to the original 'team' value ('D'). This mechanism is vital for ensuring that only the targeted values undergo the transformation, maintaining the integrity of the rest of the dataset. This particular example performs the following specific transformations within

the **team** column of the DataFrame:

Replace 'A' with 'Atlanta'

Replace 'B' with 'Boston'

Replace 'C' with 'Chicago'

## Step-by-Step Practical Example: Data Preparation

To demonstrate the application of this conditional replacement logic, let us establish a working PySpark environment and create a sample DataFrame. This DataFrame will represent data concerning basketball players, containing abbreviated team names that we intend to expand into their full city names for improved clarity and reporting. The first step involves initializing a SparkSession, which is the entry point for all Spark functionality.

We define the structure of our sample data as a list of lists, where each inner list represents a row, followed by defining the column schema. The columns include 'team' (abbreviated name), 'conference', and 'points'. This setup provides a simple yet effective scenario to test our data transformation strategy. The creation of the DataFrame is finalized using `spark.createDataFrame(data, columns)`, preparing the data structure for immediate manipulation.

The following code snippet illustrates the setup process, resulting in the display of the initial DataFrame structure, which clearly shows the abbreviated team identifiers that require replacement. This baseline view is essential for confirming the success of the subsequent transformation steps.

```
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|
+----+-----+-----+
| A| East| 11|
| A| East| 8|
| A| East| 31|
| B| West| 16|
| B| West| 6|
| C| East| 5|
| D| West| 12|
| D| West| 24|
+----+-----+-----+
```

## Executing the Value Replacement Logic

With the source DataFrame prepared, we now apply the conditional replacement logic using the `when()` and `otherwise()` functions. This transformation aims to convert the single-letter team codes ('A', 'B', 'C') into their full city names ('Atlanta', 'Boston', 'Chicago') within the 'team' column. As discussed, we must import the `when` function from `pyspark.sql.functions` before executing the transformation.

The core operation is performed using `df.withColumn()`, targeting the 'team' column. The chained `when()` clauses define the specific mappings. The use of the `otherwise(df.team)` clause is critical here, as it ensures that the team code 'D', which is not explicitly mapped, retains its original value. If the `otherwise()` clause were omitted or defined as a literal value (like `null`), all rows corresponding to 'D' would be incorrectly modified, potentially leading to data corruption or missing values.

Executing the following code block generates a new DataFrame, `df_new`, which reflects the successful application of the replacement logic. Viewing this new DataFrame confirms that the targeted values have been expanded correctly, demonstrating the efficiency and precision of conditional replacements in PySpark using this advanced syntax.

```
from pyspark.sql.functions import when
```

```
#replace multiple values in 'team' column
df_new = df.withColumn('team', when(df.team=='A', 'Atlanta')
    .when(df.team=='B', 'Boston')
    .when(df.team=='C', 'Chicago'))
    .otherwise(df.team))

#view new DataFrame
df_new.show()

+-----+-----+-----+
| team|conference|points|
+-----+-----+-----+
|Atlanta| East| 11|
|Atlanta| East| 8|
|Atlanta| East| 31|
| Boston| West| 16|
| Boston| West| 6|
|Chicago| East| 5|
| D| West| 12|
| D| West| 24|
+-----+-----+-----+
```

## Analyzing the Results and Handling Unmatched Values

Upon reviewing the output of `df_new.show()`, it is immediately apparent that the transformation has been successful. The single-letter codes 'A', 'B', and 'C' in the **team** column have been accurately replaced with 'Atlanta', 'Boston', and 'Chicago', respectively, demonstrating that the conditional logic defined by the `when()` chain executed as intended. This process confirms the reliability of using cascaded conditional statements for mapping known values to their replacements in a scalable PySpark environment.

A particularly important observation pertains to the value 'D'. Notice that we did not specify a replacement mapping for 'D' in our series of `when()` clauses. Because the value 'D' did not satisfy any of the explicit conditions (it is not 'A', 'B', or 'C'), the logic fell through to the final `otherwise(df.team)` statement. As a result, the value 'D' simply remained unchanged in the new DataFrame. This demonstrates the critical role of `otherwise()`: it provides a safety net, ensuring that data points not explicitly targeted for transformation are gracefully preserved. If we had intended to categorize 'D' (perhaps as 'Dallas'), we would have simply added another `.when(df.team=='D', 'Dallas')` before the `otherwise()` clause.

For those seeking deeper insight into the functional details, the complete documentation for the PySpark `when` function provides comprehensive details on its behavior, including its limitations and advanced usages, such as handling nullable types or complex nested conditions. Understanding this documentation is crucial for advanced PySpark development, especially when dealing with highly specific data quality or transformation requirements.

## Alternative Method: Using PySpark's `replace()` Function

While the `when()` and `otherwise()` method is generally preferred for its flexibility and similarity to SQL's `CASE WHEN` structure, PySpark does offer a more direct function for mapping replacements: the `replace()` function, available within the column operations. This function is ideally suited for scenarios where you have a simple, direct mapping of old values to new values, often expressed as a dictionary. It streamlines the process by removing the need for explicit conditional checks for each value.

The `replace()` function takes a dictionary where keys are the values to be found and replaced, and the values are the replacement values. For instance, to replace all occurrences of 1, 2, and 3 with 4, 5, and 6 respectively, a dictionary containing the mapping `{1: 4, 2: 5, 3: 6}` can be passed to the function. The function efficiently applies this mapping across the specified column. This method is often cleaner and more concise than the `when()` chain when the replacement logic is purely a one-to-one or many-to-one lookup, without complex logical dependencies.

However, a key limitation of the `replace()` function compared to the `when()` chain is its handling of values that are not present in the dictionary. Unlike the implicit preservation offered by `otherwise()`, the behavior of `replace()` might require extra steps if you need to retain values outside of the defined dictionary mapping, or if you need to handle nulls specifically. Therefore, for transformation tasks involving conditional logic, varying output types, or sophisticated defaults, the `when()` approach remains the industry standard, offering superior control and clarity in the data transformation process.

## Performance Considerations for Large Datasets

When working with truly massive datasets in a distributed environment like Apache Spark, performance considerations are paramount. Both the `when()` chain and the `replace()` function are designed to operate efficiently by leveraging Spark's internal optimizations and executing on partitioned data. They minimize data shuffling and maximize parallelism, which is superior to implementing similar logic using User Defined Functions (UDFs). UDFs often serialize and deserialize data, which introduces overhead and is generally discouraged for simple column transformations.

Between the two methods, the `replace()` function, when used with direct dictionary lookups on a

single column, can sometimes offer marginal performance advantages because it represents a direct, optimized mapping operation. However, this difference is often negligible in comparison to the flexibility gained by using the `when()` chain. The `when()` structure is compiled by the Spark Catalyst Optimizer into highly efficient physical plans, effectively achieving performance comparable to specialized SQL `CASE WHEN` expressions. Therefore, prioritizing code clarity and logical correctness, often achieved through `when()`, is typically recommended over optimizing for the slight performance edge of `replace()`, unless the dictionary mapping is exceptionally large.

Engineers should strive to minimize the number of passes over the data. Since both of these methods achieve all replacements in a single `withColumn()` operation, they are inherently efficient. The use of vectorized operations provided by `pyspark.sql.functions` ensures that data processing remains columnar, minimizing memory footprint and maximizing CPU utilization across the cluster. This adherence to native Spark functions is the single most important factor in achieving high performance in PySpark transformations.

## Conclusion: Mastering Data Transformation in PySpark

Replacing multiple values in a `DataFrame` column is a fundamental requirement for data cleansing and feature engineering in big data analytics. PySpark provides robust and scalable methods to handle this requirement. While the direct `replace()` function offers a quick solution for simple dictionary mappings, the combination of the `when()` and `otherwise()` functions, imported from `pyspark.sql.functions`, provides unparalleled control, clarity, and flexibility for complex conditional transformations.

By chaining multiple `when()` conditions and concluding with an `otherwise()` clause that preserves unmatched values, data engineers can reliably implement intricate mapping logic in a highly optimized manner. This approach ensures that data integrity is maintained, and transformation code remains readable and easily maintainable, even as the mapping requirements grow in complexity. Mastering this conditional logic is a vital skill for anyone working with data processing at scale using PySpark.

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