

# How to Split a String Column into Multiple Columns in PySpark

Authored by  
**stats writer**

January 19, 2026

## RECOMMENDED CITATION

stats writer (2026). *How to Split a String Column into Multiple Columns in PySpark*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=126655>

## 1. The Necessity of String Splitting in PySpark

Working with raw data often involves handling composite fields where multiple pieces of information are concatenated into a single PySpark DataFrame column, usually separated by a specific delimiter. Extracting these distinct components is a fundamental step in data cleaning and preparation, essential for effective analysis and modeling. In the realm of large-scale data processing powered by PySpark, the ability to efficiently manipulate these string columns is critical for maintaining performance and scalability. This process of decomposition allows data scientists and engineers to structure data appropriately, turning complex strings into manageable, granular features.

The challenge lies in applying this transformation across massive distributed datasets without resorting to slow, row-by-row operations. Fortunately, PySpark provides highly optimized, built-in functions specifically designed for vectorized operations, ensuring that transformations are executed across the cluster efficiently. By leveraging these native functions, we can effectively split single string columns into multiple derived columns, isolating meaningful data points such as geographical location, categorical identifiers, or timestamps embedded within a single string structure. This capability is paramount when dealing with semi-structured data sources like logs or scraped text where hierarchical information is flattened into delimited fields.

This guide will demonstrate the precise, robust methodology for splitting a string column within a PySpark DataFrame using the powerful built-in utilities. We will focus on the mechanics of the primary function responsible for this transformation, illustrating how to define the splitting criteria and how to correctly access the resulting elements to populate new columns. Mastering this technique ensures data preparation remains agile and compliant with the requirements of distributed computing environments.

## 2. Understanding the PySpark split() Function

The core functionality for achieving string decomposition in PySpark resides in the split() function, which is imported from the **pyspark.sql.functions** module. This function operates on a column of strings and applies a specified separation rule across all rows simultaneously. Unlike standard Python string operations, the split() function is optimized for Spark's distributed architecture, making it the preferred method for high-throughput data processing tasks.

The primary inputs required by the split() function are the target column and the delimiter. This delimiter can be a simple character, such as a hyphen (-), comma (,), or pipe (|), or it can be a complex regular expression pattern for more intricate splitting requirements. Crucially, the function does not return individual strings; instead, it returns a single column containing an array of strings for each row. This resulting array holds all the segments generated after the split operation.

To finalize the operation and extract these segments into distinct columns, we must interact with the resulting array structure. This is accomplished using the `getitem()` method, which allows indexing into the array generated by the split function. Since array indexing in computing starts at zero (**0**), the first split segment is accessed via `getitem(0)`, the second via `getitem(1)`, and so forth. By chaining these operations with the `withColumn` function, we can dynamically create as many new columns as there are segments in the split string.

### 3. Basic Syntax for Column Decomposition

To implement string splitting, we combine the `split()` function with the `withColumn()` function, which is used to add or replace columns in a [PySpark DataFrame](#). The standard process involves calling `split()` on the original column, immediately followed by `getitem()` to select the required element index, and wrapping this expression within `withColumn()` to assign a new column name.

The conceptual syntax below demonstrates how to extract two components--a 'location' and a 'name'--from an existing 'team' column, assuming the components are separated by a dash (-). This pattern is highly reusable for any string column decomposition task where the structure is consistent across rows.

You can use the following concise syntax to split a source string column into multiple derived columns within a [PySpark DataFrame](#):

```
from pyspark.sql.functions import split

#split team column using dash as delimiter
df_new = df.withColumn('location', split(df.team, '-').getitem(0))
.withColumn('name', split(df.team, '-').getitem(1))
```

In this specific example, we utilize the `split()` function twice. The first instance extracts the segment at index **0**, corresponding to the **location**, and the second extracts the segment at index **1**, corresponding to the **name**. Both operations target the original **team** column in the DataFrame, effectively creating two new, independent columns derived from the original composite string field.

### 4. Setting up the PySpark Environment and Sample Data

Before executing the splitting logic, we must ensure a [PySpark](#) session is initialized and that we have a sample [PySpark DataFrame](#) ready for transformation. For illustrative purposes, we will create a small dataset containing basketball team information where the city and team name are combined using a hyphen `delimiter`. This scenario perfectly mirrors real-world situations where compound identifiers must be broken down.

The setup involves importing **SparkSession**, creating the session object, defining our raw data list (which contains the composite string and an integer value for points), and specifying the column schema. This preparatory step is fundamental to any PySpark workflow and ensures reproducibility of the results.

The following code block demonstrates the setup required to instantiate the initial DataFrame, which we will call **df**. Notice how the **team** column holds the concatenated string we intend to separate.

## Example: Initializing the Sample DataFrame

```
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|points|
```

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

```
| Dallas-Mavs| 18|
```

```
| Brooklyn-Nets| 33|
```

```
| LA-Lakers| 12|
```

```
|Houston-Rockets| 15|
```

```
| Atlanta-Hawks| 19|
```

```
| Boston-Celtics| 24|
```

| Orlando-Magic| 28|

+-----+-----+

The output clearly shows the initial structure. We are targeting the strings within the **team** column, aiming to separate the city prefix (e.g., 'Dallas') from the team suffix (e.g., 'Mavs'). This transformation is necessary to enable filtering or aggregation based purely on the geographical location or the team nickname, tasks that would be cumbersome if the data remained concatenated.

## 5. Step-by-Step Implementation of Column Splitting

With the initial DataFrame (**df**) prepared, we can now apply the core logic using the `split()` function. The goal is to produce a new DataFrame, **df\_new**, that contains the original columns plus two newly derived columns: **location** and **name**. This process requires importing the necessary function and then executing the chained **withColumn** transformations.

For each new column, we use the following steps: first, invoke `split(df.team, '-')` to generate the array of strings using the hyphen as the separator. Second, we apply `getItem(0)` to retrieve the first element (the location) for the **location** column, and `getItem(1)` to retrieve the second element (the name) for the **name** column. Chaining these operations is highly efficient as Spark optimizes the sequence of transformations.

The following comprehensive code block executes the splitting logic and then displays the resulting DataFrame, allowing us to immediately verify the success of the transformation. Note the clean separation of the composite field into its constituent parts.

```
from pyspark.sql.functions import split
```

```
#split team column using dash as delimiter
df_new = df.withColumn('location', split(df.team, '-').getItem(0))
          .withColumn('name', split(df.team, '-').getItem(1))
```

```
#view new DataFrame
df_new.show()
```

```
+-----+-----+-----+-----+
| team|points|location| name|
+-----+-----+-----+
| Dallas-Mavs| 18| Dallas| Mavs|
| Brooklyn-Nets| 33|Brooklyn| Nets|
| LA-Lakers| 12| LA| Lakers|
```

```
|Houston-Rockets| 15| Houston|Rockets|  
| Atlanta-Hawks| 19| Atlanta| Hawks|  
| Boston-Celtics| 24| Boston|Celtics|  
| Orlando-Magic| 28| Orlando| Magic|  
+-----+-----+-----+-----+
```

## 6. Analyzing the Result and Data Transformation

Upon examining the resultant DataFrame, `df_new`, we can confirm that the string manipulation was successful. The initial `team` column, which contained values like 'Dallas-Mavs', remains in the DataFrame, but we now have two highly structured columns, `location` and `name`, populated with the discrete components of the original string. This transformation is pivotal for subsequent data operations, enabling precise filtering and aggregation based on these new features.

Specifically, the mechanism employed utilized the `split()` function to break down each composite string wherever the specified `delimiter` (the dash) appeared. This yielded a list-like structure, or an `array of strings`, for every record. We then explicitly instructed Spark which element to extract using the zero-based index provided to the `getItem()` method. The use of `getItem(0)` systematically retrieved the first part (the location), while `getItem(1)` retrieved the second part (the name).

This approach highlights the power of functional programming in `PySpark`. By defining transformations as expressions on the columns themselves, we benefit from Spark's underlying Catalyst optimizer, which plans the most efficient execution across the cluster. This is significantly more scalable and performant than attempting to iterate over rows in standard Python loops. The result is clean, structured data ready for advanced analytical tasks, such as calculating the average points scored by teams in a specific location.

## 7. Advanced Considerations and Best Practices

While the basic usage of the `split()` function is straightforward, advanced scenarios require careful attention to potential edge cases and performance considerations. One critical aspect is handling strings that may not contain the `delimiter`, or strings that contain multiple delimiters. If a row lacks the delimiter, the `split()` function returns an array containing the original string as a single element at index 0. If you attempt to access an index that does not exist (e.g., `getItem(2)` when only two segments were created), Spark will populate that new column with `null` values. Robust data pipelines often incorporate preprocessing steps to ensure consistency or utilize functions like `coalesce()` to handle expected nulls gracefully.

For scenarios involving complex separators, the `split()` function accepts a `regular expression pattern` instead of a simple character `delimiter`. This vastly expands its utility, allowing you to split

based on patterns such as sequences of whitespace, or based on specific non-alphanumeric characters, offering precise control over the tokenization process. However, using complex regular expressions can introduce computational overhead, so it is advisable to use simple character delimiters whenever possible for optimal performance in a distributed environment.

A notable optimization technique when splitting a string into many columns is to avoid redundant calculations. In the example provided, we called the **split()** function twice (once for **location** and once for **name**). For DataFrames where dozens of columns need to be extracted, it is much more efficient to perform the split operation once, storing the resulting array of strings in a temporary column, and then using getItem() repeatedly on that temporary column. Once all segments are extracted, the intermediate array column can be safely dropped. This strategy significantly reduces the amount of recomputation required by Spark.

## 8. Further PySpark Data Manipulation Techniques

Mastery of string splitting is often a stepping stone to more sophisticated data manipulation tasks within PySpark. Once strings are successfully split, data professionals frequently need to cast the resulting components to appropriate data types (e.g., converting a segment representing a timestamp or a numerical identifier from string to a native type). Other common tasks include concatenating disparate fields back together using the **concat()** function, or utilizing **regexp\_extract()** when only a specific substring needs to be isolated without performing a full split.

The official documentation for the split() function provides comprehensive details on its parameters and behavior, especially concerning different regular expression flags and handling limitations imposed by the underlying distributed file system structure. Consulting these resources is highly recommended for developers building production-level ETL pipelines.

**Note:** You can find the complete documentation for the PySpark split function online.

For readers interested in expanding their proficiency in advanced PySpark data preparation, consider exploring these related topics:

How to handle missing values using imputation methods.

Techniques for joining multiple DataFrames efficiently.

Converting data types using **cast()** functions for numerical analysis.

Implementing window functions for complex analytical calculations.