

How to Convert a Date Column to a String in PySpark

Authored by
stats writer

January 3, 2026

RECOMMENDED CITATION

stats writer (2026). *How to Convert a Date Column to a String in PySpark*.

PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=110560>

Introduction: Why Date-to-String Conversion is Essential in Data Processing

When working with large-scale data utilizing [PySpark](#), managing data types correctly is paramount for efficient processing and accurate analysis. While dates are typically stored using the Spark [DateType](#) or [TimestampType](#) for calculations, there are numerous scenarios where converting a date column into a string format becomes necessary. This is especially true when preparing data for external systems, generating human-readable reports, or interfacing with legacy applications that require specific [date format](#) representations.

The core challenge in this conversion lies not just in changing the underlying data type, but in defining the exact output format of the resulting string. Unlike simple numeric conversions, date-to-string transformation requires specifying a pattern (e.g., MM/dd/yyyy, yyyy-MM-dd) to dictate how the date components--month, day, and year--will appear in the final text output. This detailed control is critical for maintaining data integrity and ensuring consistency across various data pipelines.

Although some might instinctively look towards the generic `cast()` function, the most robust and controlled method within the [PySpark SQL functions](#) library is the use of [date_format](#). This function is specifically designed to handle the nuances of date formatting, providing developers with precision that simpler functions often lack. The following guide provides an in-depth explanation and practical example of how to leverage this powerful function to successfully convert date columns into strings in your Spark [DataFrame](#).

Understanding PySpark Data Types: DateType vs. StringType

In [PySpark](#), a [DataFrame](#) utilizes a strict schema to manage the data types of its columns. The [DateType](#) stores dates internally as the number of days since the epoch (1970-01-01), optimizing them for date arithmetic and comparisons. It does not carry any inherent formatting information; the display format seen in `df.show()` is merely a standard default chosen by Spark.

Conversely, when data is stored as a [StringType](#), it is treated purely as a sequence of characters. While this loses the mathematical properties of a date, it gains flexibility in presentation. If you need to output the date specifically as "October 30, 2023" or "30/10/23", using [date_format](#) to convert the [DateType](#) to a [StringType](#) is the required methodology.

It is important to remember the direction of conversion: if you are moving from [DateType](#) to [StringType](#), you are defining the output presentation. If you were moving from [StringType](#) to [DateType](#) (which often requires the `to_date()` function), you would be defining the input format so Spark knows how to parse the text into a valid date object. For the purpose of outputting a formatted string, we rely on the formatting capabilities built into the [date_format](#) function.

The PySpark `date_format()` Function Explained

The `date_format` function is a highly versatile tool within the `pyspark.sql.functions` module. Its purpose is explicitly to convert a date, timestamp, or string column representing a date into a string formatted according to a specified pattern. This function takes two mandatory arguments: the column containing the date value, and the string pattern defining the desired output format.

The required syntax for using this function is concise and powerful:

```
from pyspark.sql.functions import date_format
```

```
df_new = df.withColumn('date_string', date_format('date', 'MM/dd/yyyy'))
```

This example demonstrates its power: it selects the existing date column (`'date'`), applies the formatting pattern (`'MM/dd/yyyy'`), and creates a new column (`'date_string'`) populated with the resulting `StringType` values.

The formatting pattern string uses standard SQL conventions. For instance, `'MM'` represents the month number (01-12), `'dd'` represents the day of the month, and `'yyyy'` represents the four-digit year. By controlling these characters, you gain complete control over the final presentation of the date within your `DataFrame`. It is crucial to import this function explicitly from `pyspark.sql.functions` before attempting to use it.

Step-by-Step Implementation Guide

To perform this conversion effectively, the process should be broken down into three logical steps: first, importing the necessary library components; second, applying the transformation using `withColumn`; and third, verifying the result using the `dtypes` attribute. This structure ensures clean, reproducible code and reliable results within your `PySpark` environment.

The initial step involves importing the `date_format` function. This is standard practice in Spark, ensuring that the environment is ready for SQL-like manipulations. Without this import, the function call will fail, regardless of whether you are working in a notebook environment or a structured application.

The conversion itself leverages the `DataFrame` method `withColumn`. This method is non-mutating; it creates a new `DataFrame` with the specified column added or replaced. We use `df.withColumn(NewColumnName, date_format(ExistingDateColumn, FormatString))` to achieve the date-to-string transformation. If you wish to overwrite the original date column with the new string representation, you would use the original column name as the `NewColumnName`. However, best practice often dictates creating a new column to preserve the original `DataType` for

potential future computations.

Practical Example: Converting Sales Data

Let us consider a practical scenario where we have a PySpark DataFrame containing sales records. This DataFrame currently stores the date of the sale using the efficient DateType. We are tasked with generating a report where the date must be displayed in the common U.S. format (MM/dd/yyyy), necessitating the conversion to a StringType.

First, we set up our Spark session and create the initial sample data using Python's `datetime` library to ensure the 'date' column starts with the correct `DateType`:

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

```
import datetime
```

```
#define data
```

```
data = ,
```

```
,
```

```
,
```

```
,
```

```
]
```

```
#define column names
```

```
columns =
```

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

```
df = spark.createDataFrame(data, columns)
```

```
#view dataframe with full column content
```

```
df.show()
```

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

```
| date|sales|
```

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

```
|2023-10-30| 136|
```

```
|2023-11-14| 223|
```

```
|2023-11-22| 450|
```

```
|2023-11-25| 290|
```

```
|2023-12-19| 189|
```

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

Before proceeding with the conversion, it is prudent to confirm the initial data types using the `dtypes` attribute. This is a crucial validation step to ensure that the source column is indeed a `DateType`, preventing potential errors during the formatting process.

#check data type of each column

df.dtypes

As confirmed by the output, the `date` column is currently stored as a `date` type, confirming our preparation is correct.

Executing the Conversion using `date_format`

Now we apply the conversion logic. We import `date_format` and use `withColumn` to generate a new column called `date_string`. The crucial format string here is `'MM/dd/yyyy'`, which determines the textual representation of the date values in the new column.

The complete execution block demonstrates the transformation and immediate visualization of the resulting `DataFrame`, showing both the original date and the newly formatted string side-by-side:

from pyspark.sql.functions import date_format

```
#create new column that converts dates to strings
df_new = df.withColumn('date_string', date_format('date', 'MM/dd/yyyy'))
```

```
#view new DataFrame
```

```
df_new.show()
```

```
+-----+-----+-----+
| date|sales|date_string|
+-----+-----+-----+
|2023-10-30| 136| 10/30/2023|
|2023-11-14| 223| 11/14/2023|
|2023-11-22| 450| 11/22/2023|
|2023-11-25| 290| 11/25/2023|
|2023-12-19| 189| 12/19/2023|
+-----+-----+-----+
```

As observed in the output, the `date_string` column now holds the dates in the desired `'MM/dd/yyyy'` format. This transformation is achieved efficiently across all rows of the `DataFrame`, leveraging Spark's optimized distributed processing capabilities. This new column is ready for

reporting or any operations that specifically require a string representation of the date.

Verification: Checking the New Data Type

The final and equally important step is to confirm that the conversion was successful and that the new column, `date_string`, is indeed recognized by PySpark as a `StringType`. This verification step prevents downstream errors in ETL processes where strict data type adherence is enforced. We use `df_new.dtypes` on our resultant DataFrame.

The output confirms the success of the operation:

```
#check data type of each column  
df_new.dtypes
```

We can clearly see that the original `date` column retains its `date` type, while the newly created `date_string` column has the expected `string` type. This validates that the `date_format` function successfully performed the transformation according to the specified format pattern, resulting in a compliant `StringType` output.

Customizing Date Format Patterns

The true flexibility of using `date_format` lies in the extensive range of format patterns supported. While we used `'MM/dd/yyyy'` in the example, which is common in the U.S., developers frequently need to conform to other global standards or specific reporting requirements. Understanding the common pattern symbols is crucial for effective customization.

For instance, if we needed the `date_format` to adhere to the ISO 8601 standard, the pattern would be `'yyyy-MM-dd'`. If the requirement was a more verbose format, such as including the abbreviated month name (e.g., Nov 14, 2023), the appropriate pattern would be `'MMM dd, yyyy'`. Other useful symbols include `'D'` for the day of the year, `'W'` for the week of the year, and `'E'` for the day name.

It is highly recommended that users consult the official Spark documentation regarding date and time patterns, as compatibility between different database standards (like Java `SimpleDateFormat` or SQL format models) can vary slightly. Choosing the correct pattern ensures that the resulting `StringType` is exactly what the downstream application expects, preventing parsing failures or misinterpretation of data.

Conclusion: Leveraging `date_format` for Robust PySpark ETL

Converting a column from `DateType` to `StringType` in `PySpark` is a fundamental operation in many Extract, Transform, Load (ETL) pipelines. By utilizing the dedicated `date_format` function, developers ensure not only the successful type conversion but also precise control over the resulting date presentation.

The method described--importing `date_format`, applying it via `withColumn` with a custom format string, and verifying the result with `dtypes`--is the industry standard for achieving this transformation cleanly and efficiently within a distributed environment like Apache Spark. Remember that the choice of the `date format` string (e.g., `MM/dd/yyyy`, `yyyy-MM-dd`) is critical and must align with the target system requirements.

By mastering this technique, you can confidently prepare your Spark `DataFrame` data for any destination system, whether that involves generating concise reports, exporting data to flat files, or integrating with diverse APIs that mandate specific string representations of temporal information.

You can use the following syntax to convert a column from a date to a string in `PySpark`:

```
from pyspark.sql.functions import date_format

df_new = df.withColumn('date_string', date_format('date', 'MM/dd/yyyy'))
```

This particular example converts the dates in the `date` column to strings in a new column called `date_string`, using `MM/dd/yyyy` as the `date format`.

The following example shows how to use this syntax in practice.

Example: How to Convert Column from Date to String in PySpark

Suppose we have the following `PySpark DataFrame` that contains information about sales made on various days for some company:

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

```
import datetime
```

```
#define data
```

```
data = ,
```

```
,
```

```
,
```

```
,  
]  
  
#define column names  
columns =  
  
#create dataframe using data and column names  
df = spark.createDataFrame(data, columns)  
  
#view dataframe with full column content  
df.show()
```

```
+-----+-----+  
| date|sales|  
+-----+-----+  
|2023-10-30| 136|  
|2023-11-14| 223|  
|2023-11-22| 450|  
|2023-11-25| 290|  
|2023-12-19| 189|  
+-----+-----+
```

We can use the **dtypes** function to check the data type of each column in the [DataFrame](#):

```
#check data type of each column  
df.dtypes
```

We confirm that the **date** column currently has a data type of [DateType](#).

To convert this column from a date to a string, we use the following syntax leveraging the [date_format](#) function:

```
from pyspark.sql.functions import date_format  
  
#create new column that converts dates to strings  
df_new = df.withColumn('date_string', date_format('date', 'MM/dd/yyyy'))  
  
#view new DataFrame  
df_new.show()
```

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

```
| date|sales|date_string|
+-----+-----+-----+
|2023-10-30| 136| 10/30/2023|
|2023-11-14| 223| 11/14/2023|
|2023-11-22| 450| 11/22/2023|
|2023-11-25| 290| 11/25/2023|
|2023-12-19| 189| 12/19/2023|
+-----+-----+-----+
```

We use the **dtypes** function once again to view the data types of each column in the resultant DataFrame:

```
#check data type of each column
df_new.dtypes
```

We can see that the **date_string** column now holds a data type of string, completing the conversion process successfully.

We have successfully created a string column from a date column while applying specific formatting.

Note: We used **MM/dd/yyyy** as the date format within the date_format function but feel free to use whatever date format you'd like based on your project requirements.

Related PySpark Tutorials

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