

# Convert String to Timestamp in PySpark (With Example)

Authored by  
**stats writer**

November 16, 2025

## RECOMMENDED CITATION

stats writer (2025). *Convert String to Timestamp in PySpark (With Example)*.  
PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=92506>

## Introduction to Date/Time Handling in PySpark

Processing time-series data or log files often requires robust handling of date and time information. In the big data world, PySpark DataFrames are the foundational structure for data manipulation. However, data frequently enters the system with dates and times stored as simple String type columns, which limits computational capabilities and performance.

A String representation, while human-readable, lacks the inherent ordering, arithmetic capabilities, and time zone awareness necessary for complex analytical tasks. To unlock the full potential of time-based analysis--such as calculating durations, filtering by time windows, or aggregating data across specific temporal boundaries--we must convert these String fields into a native time format, specifically the Timestamp data type.

This conversion is critical for maintaining data quality and ensuring that Spark's optimized internal mechanisms for date and time handling are fully utilized. This comprehensive guide details the precise methodology for performing this transformation using built-in PySpark SQL functions, ensuring your data is primed for high-performance time-series analysis.

### The Necessity of Converting String to Timestamp

Why is it so vital to switch from the generic String format to the specific Timestamp format? The primary reason lies in how PySpark DataFrames manage different data types. When a column is stored as a String, Spark treats it as raw text; operations like subtraction, comparison based on chronological order, or extraction of components (like hour or minute) become cumbersome, requiring complex text parsing or regular expressions which are computationally expensive.

The native Timestamp type, conversely, stores the value as the number of milliseconds since the epoch (January 1, 1970 UTC). This numerical representation allows for extremely fast and accurate time calculations, which is crucial in a distributed processing environment like Spark. Furthermore, PySpark provides a rich library of specialized functions optimized specifically for the Timestamp type, such as `date_add`, `datediff`, and `window`, which are inaccessible or inefficiently used when the data remains in String format.

In essence, converting the column ensures that your PySpark DataFrame is ready for deep analytical processing, adhering to best practices for high-performance data types in a distributed computing environment. This foundational step is non-negotiable for anyone performing serious time-series data engineering.

### Understanding the PySpark to\_timestamp Function Syntax

The transformation from String to Timestamp is accomplished using the robust built-in function,

`F.to_timestamp`, which resides within the `pyspark.sql.functions` module. This function requires two primary arguments: the name of the source column (the `String` field) and the precise format pattern of that string.

The general syntax structure is straightforward. We utilize the `DataFrame` method `withColumn` to either create a new column or overwrite an existing one, applying the `F.to_timestamp` function as the transformation logic. The function parses the string based on the provided pattern and returns a column of `Timestamp` type.

The format pattern is crucial. It must exactly match the structure of the date and time elements in your source `String` column. Standard codes include `YYYY` for the four-digit year, `MM` for the month, `dd` for the day, `HH` for the hour (24-hour clock), `mm` for the minute, and `ss` for the second. Any deviation in the pattern (e.g., mixing `MM` for month with `mm` for minute) or incorrect placement of separators will lead to failed parsing and result in null values in the target column.

The following snippet demonstrates the standard application of this function to convert a column named `ts` into a new column named `ts_new`:

```
from pyspark.sql import functions as F
```

```
df = df.withColumn('ts_new', F.to_timestamp('ts', 'yyyy-MM-dd HH:mm:ss'))
```

This operation ensures that all subsequent analytical functions treat `ts_new` not as text, but as a chronologically ordered time value.

## Step-by-Step Example: Initializing the PySpark Environment and Data

To illustrate the conversion process effectively, we will construct a small, representative `PySpark DataFrame` containing simulated sales data. This data set includes a time column, `ts` (representing the transaction timestamp), which is initially stored in the suboptimal `data type` of `string`.

Our scenario involves tracking sales transactions where the recording time is captured as a combined date and time `String`. Successful conversion to a proper `Timestamp` will enable us to accurately analyze sales trends over time, calculate daily or hourly averages, or identify peak sales periods with accuracy. This initial step mimics the ingestion process of raw data into a Spark environment.

The following code snippet initializes a `SparkSession` and defines our sample data. This is the crucial starting point for any PySpark operation:

```
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()
```

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

```
| ts|sales|
```

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

```
|2023-01-15 04:14:22| 225|
```

```
|2023-02-24 10:55:01| 260|
```

```
|2023-07-14 18:34:59| 413|
```

```
|2023-10-30 22:20:05| 368|
```

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

The `df.show()` output confirms that our initial data structure is loaded correctly, presenting four sales records associated with specific moments in time, all currently held within the `ts` column as visible text strings.

## Verifying Initial Data Types (The Challenge)

Before proceeding with the conversion, it is essential to explicitly confirm the existing data types defined by Spark during the DataFrame creation process. PySpark often infers data types based on the input data, and in the case of complex date/time characters embedded within quotation marks, it typically defaults to the generic String type unless an explicit schema is provided.

We inspect the schema using the `df.dtypes` attribute, which returns a list of tuples detailing the column name and its associated type. This step explicitly demonstrates the necessity of the subsequent type conversion operation, highlighting why direct time calculations cannot yet be performed on the `ts` column.

## #check data type of each column

### df.dtypes

As clearly demonstrated in the output, while the `sales` column is correctly identified as a numerical `bigint`, the crucial time column, `ts`, is categorized as a `string`. This confirmation verifies that the data is not yet optimized for temporal analysis, setting the stage for the proper application of the `F.to_timestamp` function.

## Implementing the Conversion using `to_timestamp`

The core solution involves importing the necessary functions module from `pyspark.sql` and using the `withColumn` method to append the converted data as a new column. This approach is preferred as it preserves the original `String` field (`ts`) while allowing us to introduce and analyze the new `Timestamp` field (`ts_new`) side-by-side.

The key to success here is the format string: `'YYYY-MM-dd HH:mm:ss'`. This mask must precisely map the structure of the input strings, including the year (`YYYY`), month (`MM`), day (`dd`), the space separator, and the time components (`HH:mm:ss`). Since our input strings follow this exact structure, the parsing will be clean and accurate.

Executing the following code transforms the data and generates the new, time-optimized column within the `DataFrame`:

### from pyspark.sql import functions as F

```
#convert 'ts' column from string to timestamp
df = df.withColumn('ts_new', F.to_timestamp('ts', 'yyyy-MM-dd HH:mm:ss'))

#view updated DataFrame
df.show()
```

```
+-----+-----+-----+
| ts|sales| ts_new|
+-----+-----+-----+
|2023-01-15 04:14:22| 225|2023-01-15 04:14:22|
|2023-02-24 10:55:01| 260|2023-02-24 10:55:01|
|2023-07-14 18:34:59| 413|2023-07-14 18:34:59|
|2023-10-30 22:20:05| 368|2023-10-30 22:20:05|
+-----+-----+-----+
```

Upon visual inspection, the values in the new column, `ts_new`, look identical to the original `ts` column. However, it is imperative to remember that their underlying structure is now optimized for numerical and temporal operations, which is the functional difference we seek.

## Validating the Transformed DataFrame

The final and most critical step is to confirm that the conversion was successful by re-examining the schema of the modified PySpark DataFrame. We utilize the `df.dtypes` attribute once more to verify the internal representation of the newly created `ts_new` column.

If the `F.to_timestamp` function executed correctly with a matching format string, the output should clearly designate `ts_new` as a timestamp data type, confirming its readiness for advanced time-series operations. Any other type, or a `null` value for a non-null input string, would indicate a formatting error.

### #check data type of each column

#### `df.dtypes`

The resulting schema confirms the successful conversion: `ts_new` is now properly recognized as a `timestamp` type. We have successfully completed the objective of transitioning the temporal data from an inefficient text format to Spark's optimized native format.

## Best Practices and Common Pitfalls

While the `F.to_timestamp` function is robust, practitioners should be aware of several best practices and common pitfalls. The most frequent error encountered is a mismatch between the input String format and the provided format mask. Developers must pay close attention to case sensitivity (e.g., `MM` for month vs. `mm` for minute) and the presence of separators. If the input format contains time zone information (e.g., `'Z'` or `'+00:00'`), this must be included in the format string (e.g., `'yyyy-MM-dd HH:mm:ssXXX'`).

Another critical consideration is handling null values. If `F.to_timestamp` fails to parse a String (due to incorrect format or corrupted data), it defaults the output Timestamp value to `null`. It is often necessary to preemptively cleanse or filter the input strings to prevent widespread null generation, or use conditional logic if multiple date string formats exist within the same column.

Finally, remember that PySpark's Timestamp type stores a point in time relative to UTC. When converting from a String that lacks explicit time zone information, Spark applies the session's default local timezone before converting the time to UTC for storage. Always confirm your Spark session configuration regarding time zones to avoid subtle data inaccuracies in cross-regional

analyses.

ARABPSYCHOLOGY.COM