

# How to Convert a Wide DataFrame to Long Format in PySpark

Authored by  
**stats writer**

January 2, 2026

## RECOMMENDED CITATION

stats writer (2026). *How to Convert a Wide DataFrame to Long Format in PySpark*.  
PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=110490>

Transforming data structures is a fundamental requirement in large-scale data processing, especially when working with big data tools like [PySpark DataFrames](#). Reshaping a [DataFrame](#) from a [wide format](#) to a [long format](#) (often called unpivoting) is a common operation necessary for analytical tasks, machine learning preparation, or visualization. Historically, achieving this transformation in [PySpark](#) required complex combinations of functions like `selectExpr`, `groupBy`, `pivot`, and `agg`. However, modern [PySpark](#) versions, particularly those supporting the powerful [melt function](#), simplify this process dramatically.

## Introduction to DataFrame Reshaping

Data [reshaping](#), or data transformation, involves altering the arrangement of rows and columns in a dataset while preserving the underlying information. This process is crucial because different analytic tasks demand specific data layouts. For instance, statistical modeling libraries often require variables (features) to be arranged vertically, which corresponds to the [long format](#). Conversely, data entry or human-readable summary tables are often stored in the [wide format](#), where each observation might span a single row, and related measures occupy separate columns.

When converting from [wide to long](#), we essentially take multiple value columns and consolidate them into two new columns: a variable column (which holds the original column names) and a value column (which holds the corresponding data entries). The remaining identifier columns, known as ID variables, remain unchanged and are repeated across the new, expanded rows. Understanding this conceptual shift is the first step toward mastering data manipulation in large-scale environments like [PySpark](#).

The original approach relying on `selectExpr`, `groupBy`, `pivot`, and `agg` was powerful but verbose and often non-intuitive for users familiar with reshaping tools in R or pandas. The introduction of the dedicated [melt function](#) standardizes and simplifies this common transformation, allowing data engineers and analysts to write cleaner, more maintainable code when preparing data for subsequent analysis.

## Understanding Wide vs. Long Data Formats

The distinction between the [wide data format](#) and the [long data format](#) is central to effective data management. In the [wide format](#), a single observational unit occupies one row, and attributes measured over different time points or categories are spread across distinct columns. For example, if measuring the performance (points) of three different positions (Guard, Forward, Center) for a single team, the wide format would have one row per team and three separate columns for the points scored by each position.

Conversely, the [long format](#) consolidates those attributes into fewer columns but increases the number of rows. Each measurement (e.g., Guard points, Forward points, Center points) becomes

its own unique row. This format typically consists of three column types: one or more ID columns (identifying the unit, like 'team'), a variable column (identifying the type of measurement, like 'position'), and a value column (containing the measurement itself, like 'points'). The long format is often preferred for analysis because it adheres to the principles of tidy data, making it easier to filter, group, and calculate summary statistics across the measurement types.

Choosing the correct format depends heavily on the destination tool or specific analytical need. While the wide format is excellent for visual comparison of metrics side-by-side, the long format is indispensable for time-series analysis, plotting multiple metrics using visualization tools that require a single grouping column, and interfacing with statistical packages that expect variables in a stackable, row-wise structure.

## Why Reshaping DataFrames is Essential for Analysis

Data analysis workflows frequently necessitate data reshaping to meet the input requirements of advanced tools. For instance, many supervised machine learning algorithms, which are often executed using libraries compatible with PySpark, expect features to be vectorized and ordered in a way that is naturally facilitated by the long format. If a dataset contains dozens of columns representing repeated measures (e.g., monthly sales figures for three years), converting these 36 columns into a single 'Month' variable column and a single 'Sales Value' column simplifies feature engineering and model training.

Furthermore, data reshaping improves database efficiency and consistency. When data is stored in the long format, it adheres closer to normalization principles, reducing redundancy. In the context of big data processing on systems like Spark, minimizing redundant information and structuring data optimally for parallel processing can yield significant performance benefits. The shift from a row containing dozens of loosely related attributes to multiple rows linked by a key dramatically enhances the flexibility of querying and filtering.

If analysts attempted to perform complex filtering or aggregations on the wide format, they would often need repetitive logic across many columns (e.g., calculating the average across Guard, Forward, and Center columns separately). The long format allows a single aggregation function applied to the 'value' column, grouped by the newly created 'variable' column, drastically simplifying complex analytical queries and ensuring statistical computations are applied uniformly across all categories.

## PySpark's `melt` Function: The Modern Approach

The melt function in PySpark is the dedicated method for performing unpivoting (wide-to-long) operations. It is designed to be highly intuitive, mirroring similar functionality found in pandas, making the transition for data scientists already familiar with that ecosystem seamless. The primary

goal of `df.melt()` is to take a DataFrame where some columns represent variable names and reshape it so that those variable names become values in a single column, and their corresponding data becomes values in a second column.

The efficiency of `melt` ensures that even when dealing with extremely large datasets managed by Spark, the reshaping process remains performant. By abstracting the complex SQL-like operations previously needed (involving joins, unions, and `stack` expressions), `melt` significantly improves development speed and reduces the likelihood of coding errors associated with manual column handling. Its introduction marks a major step forward in making PySpark accessible for complex data preparation tasks.

The basic structure of the `melt function` requires specifying which columns are identifiers (to keep static) and which columns contain the values to be unpivoted. The example below illustrates the foundational syntax used to initiate this transformation:

You can use the **melt function** with the following basic syntax to convert a PySpark DataFrame from a wide format to a long format:

```
df_long = df.melt(ids=, values=,  
variableColumnName='position',  
valueColumnName='points')
```

This particular example converts a wide DataFrame named `df` to a long DataFrame named `df_long`.

The following section demonstrates how to use this syntax in a complete, practical coding example.

## Detailed Syntax and Parameters of `df.melt()`

The `df.melt()` method accepts several key parameters that control how the DataFrame is reshaped. Proper configuration of these parameters is vital for ensuring the output structure aligns with analytic requirements. These parameters allow the user precise control over which columns are maintained as identifiers and how the new variable and value columns are named.

The four primary parameters of the `melt function` are:

**ids:** This parameter takes a list of column names that should remain as identifier variables. These columns are duplicated across the new, expanded rows and form the key structure of the long format output. If omitted, all non-value columns are treated as ID columns.

**values:** This parameter takes a list of column names containing the data values that need to be unpivoted. These are the columns that will be collapsed into a single 'value' column.

**variableColumnName:** A string specifying the name of the new column that will store the names of the original columns specified in the `values` list (e.g., 'position' in the example).

**valueColumnName:** A string specifying the name of the new column that will store the actual data content from the original columns specified in the `values` list (e.g., 'points' in the example).

The flexibility offered by naming the output columns using `variableColumnName` and `valueColumnName` is particularly useful for improving code readability and ensuring that the final long format DataFrame is immediately meaningful within the context of the business problem. Using descriptive names like `position` and `points` instead of generic defaults (like `variable` and `value`) enhances data governance and interpretability downstream.

## Setting Up the PySpark Environment

Before executing the reshape operation, it is necessary to initialize a Spark session and create the initial wide format DataFrame. This setup ensures that the environment is configured correctly and provides a concrete dataset to work with. Suppose we are working with basketball statistics, where the performance data for two teams (A and B) is organized by player position across columns. This is the definition of a wide data format.

The following code snippet demonstrates the necessary steps to import the required classes, instantiate the `SparkSession`, define the data, and create the initial DataFrame named `df`. Notice how the initial output shows the data spread across four columns: `team` (the ID variable), and `Guard`, `Forward`, and `Center` (the value variables).

```
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|Guard|Forward|Center|
+---+---+-----+-----+
| A| 22| 34| 17|
| B| 25| 10| 12|
+---+---+-----+-----+
```

The resulting `DataFrame` `df` is clearly in the wide format. Each team (A and B) occupies a single row, and the point totals for each position are spread horizontally. Our goal is to consolidate the position data (Guard, Forward, Center) into rows, making 'position' a new grouping variable and 'points' the corresponding value.

### Practical Example: Reshaping a Wide DataFrame

Now that the environment is set up and the initial wide DataFrame is prepared, we apply the `melt` function to perform the required data reshaping. We must carefully define the `ids` (identifier columns) and `values` (columns to unpivot).

In this scenario, 'team' is the unique identifier we want to retain across rows, so it is passed to `ids`. The columns 'Guard', 'Forward', and 'Center' contain the values we wish to stack, so they are passed to `values`. Finally, we assign explicit names to the new columns: `position` for the variable name and `points` for the data value. This precise control over column mapping results in the desired long format structure, as shown below:

```
#create long DataFrame
df_long = df.melt(ids=, values=,
variableColumnName='position',
valueColumnName='points')
```

```
#view long DataFrame
df_long.show()
```

```
+---+-----+-----+
|team|position|points|
+---+-----+-----+
| A| Guard| 22|
| A| Forward| 34|
| A| Center| 17|
| B| Guard| 25|
```

```
| B| Forward| 10|
| B| Center| 12|
+---+-----+-----+
```

## Analyzing the Reshaped (Long) Output

The resulting `DataFrame`, `df_long`, clearly demonstrates the successful transformation into the long format. Where the original data contained two rows, the reshaped data now contains six rows (two original rows multiplied by three value columns). This increase in row count is characteristic of unpivoting operations, as each individual measurement becomes its own record.

Specifically, the **team** identifier is now repeated along the rows, ensuring that each data point remains correctly associated with its source team. The original column headers ('Guard', 'Forward', 'Center') have been collapsed and their names are now used as values in the second column, which we explicitly named **position** using the `variableColumnName` argument. Finally, the actual point totals (22, 34, 17, etc.) are consolidated into the third column, which we named **points** using the `valueColumnName` argument.

This structure is highly optimized for analytical aggregation. For example, calculating the total points scored by all Guards across all teams is now a simple `groupBy('position').agg(sum('points'))` operation. In the original wide format, this operation would have been trivial for this small dataset but far more complex and prone to errors if dozens of position columns existed. The long format is thus superior for scalability and statistical processing.

## Legacy Methods: Using `stack` or `pivot`

While the melt function is the recommended modern approach for wide-to-long format transformations, earlier versions of PySpark or complex scenarios sometimes required the manual construction of the unpivot using SQL expressions or the `stack` function. The initial paragraph of this article referenced a combination of `selectExpr`, `groupBy`, `pivot`, and `agg`, which is typically used for the opposite transformation (long-to-wide), but variations involving `stack` within `selectExpr` were used for unpivoting.

The `stack` function, used within a `selectExpr` clause, manually defines the pairs of variable names and corresponding values to be unpivoted. This method is highly verbose and demands meticulous definition of column indices and names, making it brittle to changes in the source schema. For instance, unpivoting three columns requires a `stack(3, 'Guard', Guard, 'Forward', Forward, 'Center', Center)` statement. This approach lacks the clarity and ease of maintenance provided by the parameter-driven melt function, reinforcing why `melt` should be prioritized for standard wide-to-long format data reshaping.

In summary, while understanding the underlying mechanisms that necessitate data reshaping remains crucial, the `df.melt()` function provides an elegant, robust, and industry-standard solution for handling wide-to-long transformations in contemporary PySpark workflows.

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

ARABPSYCHOLOGY.COM