

How to Find the Minimum Value Across Columns in PySpark

Authored by
stats writer

February 9, 2026

RECOMMENDED CITATION

stats writer (2026). *How to Find the Minimum Value Across Columns in PySpark*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=129853>

The process of calculating the minimum value across multiple columns in [PySpark](#) is a fundamental requirement in advanced [data analysis](#) workflows, particularly when dealing with wide datasets where row-level aggregation is necessary. Unlike standard column aggregations (like finding the overall minimum score in 'game1'), cross-column aggregation involves finding the smallest value present for a specific record (or row) among a defined subset of columns. [PySpark](#) provides highly optimized native functions within the [pyspark.sql.functions](#) module to handle this task efficiently, bypassing the need for slower User Defined Functions (UDFs) that often introduce serialization overhead. The primary function employed for this purpose is **least**, which is specifically designed to compare scalar values across multiple columns and return the minimum among them. Understanding how to correctly implement this function allows for powerful data manipulation and feature engineering directly within the distributed computing framework.

PySpark: Calculate Minimum Value Across Columns

Introduction to Cross-Column Aggregation in PySpark

When working with large-scale datasets in [PySpark](#), data often represents repeated measurements or scores across different variables, structured horizontally across columns. A common analytical requirement is to derive a new feature that summarizes the performance or measurement across these variables at the row level. For instance, if a dataset tracks sales across three different regions (Region A, Region B, Region C), finding the minimum sales figure for each individual product across these three regions requires a cross-column minimum calculation. This operation is fundamentally different from a standard SQL aggregation (like `MIN()` applied to a single column), which aggregates vertically over all rows. By using optimized PySpark functions, we ensure that these calculations are performed in parallel across the cluster, maintaining the efficiency required for big data processing.

The core philosophy behind efficient operations in PySpark involves utilizing built-in functions whenever possible, as they are executed using the optimized Catalyst optimizer and run on the JVM, leading to significant performance gains over Python-based logic. The native **least** function is the canonical solution for finding the minimum value among several specified columns in a [PySpark DataFrame](#). While it is technically possible to achieve this using complex logic involving converting columns to arrays or implementing a User Defined Function (UDF), these methods introduce unnecessary complexity and are substantially less performant than relying on the specialized functions provided by the framework, making the **least** function the preferred industry standard approach.

It is critical to understand the constraints and capabilities of the cross-column minimum calculation. First, the **least** function is primarily designed for comparing numerical columns. While it can handle certain data types like strings (performing comparisons based on lexicographical order), it is most

frequently and reliably used with integer, float, or double data types. Secondly, the behavior concerning null values must be considered carefully. PySpark's **least** function is designed to ignore nulls during comparison; however, if all input columns for a given row are null, the resulting minimum value will also be null. This behavior is crucial for ensuring accurate results during complex feature engineering tasks.

The Power of the least Function

The `least` function, imported from the `pyspark.sql.functions` module, provides a straightforward, highly optimized way to achieve row-wise minimum calculation. Unlike the general `min` function which finds the minimum value across all rows within a single column, **least** operates horizontally, accepting multiple column inputs and returning the smallest value among those inputs for the current row. This distinction is vital when performing statistical summaries or data quality checks that require comparing data points within the same record.

When using **least**, the inputs must be column references--either string names corresponding to columns in the `DataFrame` or column objects created using `col()` or other PySpark functions. The function can take an arbitrary number of column arguments, making it highly flexible for datasets with varying numbers of input features. For instance, if you have thirty different measurements, you can pass all thirty column names directly into the **least** function call. The underlying Spark engine then executes this comparison efficiently in a distributed manner, minimizing data shuffling and maximizing throughput.

Beyond simple numerical comparison, the effectiveness of the **least** function lies in its inherent compatibility with Spark's execution plan. Because it is a native function, the Catalyst optimizer can efficiently translate it into optimized byte code (often Scala or Java), meaning it avoids the costly Python interpreter execution overhead. This makes it suitable not just for small demonstration examples, but for production environments handling petabytes of data. This architectural advantage is the primary reason why developers are strongly advised to favor built-in functions like **least** over custom Python implementations.

Syntax and Practical Implementation using withColumn

To integrate the minimum calculation into an existing `PySpark DataFrame`, the standard practice is to use the `withColumn` transformation. The **withColumn** method is non-mutating, meaning it returns a new `DataFrame` with the specified column added or replaced, ensuring data integrity of the original structure. The general syntax involves importing the function and then applying it within the **withColumn** call, assigning the result of **least** to a new column name.

The required steps are minimal and highly declarative. First, ensure the necessary functions are imported from `pyspark.sql.functions`. Second, call the **withColumn** method on your

DataFrame, providing the name for the new resulting column (e.g., 'min_score') as the first argument. Third, define the column expression using the **least** function, passing the names of all target columns as positional arguments. This structure clearly defines the data transformation logic, making the code readable and maintainable, a key characteristic of good PySpark development practices.

The following syntax demonstrates the standard usage pattern. This particular example creates a new column called **min** that contains the minimum of values across the **game1**, **game2** and **game3** columns in the DataFrame:

```
from pyspark.sql.functions import least
```

```
#find minimum value across columns 'game1', 'game2', and 'game3'  
df_new = df.withColumn('min', least('game1', 'game2', 'game3'))
```

This implementation immediately benefits from PySpark's lazy execution model. The calculation is not performed until an action (like `.show()` or `.collect()`) is called, allowing the Catalyst optimizer to analyze the entire execution plan and apply optimizations before the distributed computation begins. This ensures that the minimum calculation is efficiently integrated into the overall workflow, minimizing resource usage and execution time.

Detailed Example: Calculating Player Minimum Scores

To illustrate this functionality in a practical context, consider a scenario involving sports data analysis where we track the performance of basketball teams across several games. We want to identify the minimum score achieved by each team during the monitored period. This minimum score often serves as an indicator of consistency or minimum guaranteed performance, which is valuable for strategic evaluation.

We begin by setting up the necessary Spark environment and defining our dataset. This example uses a simple list of tuples to represent the data, which is then converted into a PySpark DataFrame with predefined column names, ensuring correct data types for the numerical game scores. The initial DataFrame creation phase is crucial for ensuring that subsequent operations, such as the minimum calculation, operate on well-structured and schema-validated data.

```
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|game1|game2|game3|
+-----+-----+-----+
| Mavs| 25| 11| 10|
| Nets| 22| 8| 14|
| Hawks| 14| 22| 10|
| Kings| 30| 22| 35|
| Bulls| 15| 14| 12|
| Blazers| 10| 14| 18|
+-----+-----+-----+-----+
```

Once the DataFrame is established, we proceed to calculate the minimum score across the 'game1', 'game2', and 'game3' columns using the **least** function in conjunction with withColumn. This operation adds a new column named **min** to the DataFrame, providing a row-level aggregation result. Observing the output confirms that for each team, the new **min** column accurately reflects the lowest value found among the three game columns, demonstrating the function's effectiveness.

```
from pyspark.sql.functions import least
```

```
#find minimum value across columns 'game1', 'game2', and 'game3'
df_new = df.withColumn('min', least('game1', 'game2', 'game3'))

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

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

```
| team|game1|game2|game3|min|
+-----+-----+-----+-----+----+
| Mavs| 25| 11| 10| 10|
| Nets| 22| 8| 14| 8|
| Hawks| 14| 22| 10| 10|
| Kings| 30| 22| 35| 22|
| Bulls| 15| 14| 12| 12|
| Blazers| 10| 14| 18| 10|
+-----+-----+-----+-----+----+
```

The resulting DataFrame clearly shows the computed minimums:

The minimum of points for the **Mavs** player is **10** (found in game3).

The minimum of points for the **Nets** player is **8** (found in game2).

The minimum of points for the **Hawks** player is **10** (found in game3).

The minimum of points for the **Kings** player is **22** (found in game2).

The minimum of points for the **Bulls** player is **12** (found in game3).

The minimum of points for the **Blazers** player is **10** (found in game1).

This output confirms that the `least` function successfully achieved the row-wise minimum calculation, adding valuable derived information to the dataset that can be used for further analysis, such as calculating performance variance or identifying outlier low scores.

Handling Edge Cases: Null Values and Data Types

A critical consideration when performing cross-column aggregations is how `PySpark` handles data quality issues, specifically null values and inappropriate data types. The `least` function is specifically designed to operate robustly in the presence of missing data, but its behavior must be correctly anticipated. Generally, `least` follows the principle that if any non-null value exists among the inputs, the function will successfully return the minimum of the non-null values, effectively excluding the nulls from the comparison set.

However, the exception to this rule occurs when all the input columns for a given row are null. In this specific scenario, since there are no numerical values to compare, the result of the `least` function for that row will also be `null`. This differs from some SQL implementations where complex coalescing or filtering might be needed. `PySpark`'s default behavior simplifies data cleansing steps, as rows entirely lacking data in the aggregated columns are flagged clearly with a null result in the new minimum column. Developers relying on this function should always verify the underlying schema, ensuring that the target columns are correctly cast to numerical types (`IntegerType`, `DoubleType`, etc.) before applying `least`, as comparing strings using `least` yields results based on

lexicographical order, which may not align with intended numerical minimums.

If non-numeric data were included in the calculation without conversion, Spark would attempt to use the lexicographical order (alphabetical) for comparison, which is almost always incorrect for finding the true numerical minimum. For example, comparing '10' and '5' lexicographically results in '10' being smaller because it starts with '1', whereas numerically '5' is smaller. Therefore, rigorous schema validation and type casting (e.g., using `.cast("int")` or `.cast("double")`) should precede any application of the **least** function to ensure that all target columns are standardized numerical formats, guaranteeing mathematically sound results from the cross-column minimum operation.

Scaling and Dynamic Column Selection

While the manual input of column names (e.g., `least('game1', 'game2', 'game3')`) works perfectly for small, static schemas, real-world big data scenarios often involve dozens or even hundreds of columns that require aggregation, or the list of columns might need to be dynamically determined based on filtering or user input. Manually listing hundreds of column names within the **least** function call is cumbersome and error-prone. To address this issue at scale, PySpark allows for dynamic column selection using lists and the asterisk (*) unpacking operator.

The solution involves first generating a Python list containing the names of all the columns intended for the minimum calculation. This list can be generated programmatically by inspecting the `DataFrame` schema and filtering column names based on a pattern (e.g., all columns starting with 'score_'). Once the list is compiled, it is passed to the **least** function using the * operator, which unpacks the list elements as individual positional arguments required by the function. This method drastically improves code maintainability and scalability when dealing with wide datasets.

For example, if the desired column list is stored in a variable `score_cols = []`, the dynamic application would look like this: `df.withColumn('min_score', least(*score_cols))`. This powerful idiom ensures that the code remains concise regardless of the number of columns involved, facilitating easier integration into automated ETL pipelines. Furthermore, this dynamic approach is fully compatible with the performance optimizations of the native `pyspark.sql.functions`, ensuring that scalability does not come at the cost of execution efficiency.

Summary and Performance Benefits

Calculating the minimum value across columns in `PySpark` is most effectively achieved using the built-in **least** function combined with the `withColumn` transformation. This technique provides a row-wise minimum calculation that is essential for many statistical and feature engineering tasks in big data environments. Key advantages include its inherent optimization by the Spark Catalyst engine, its ability to handle multiple column inputs dynamically, and its robust handling of null

values by ignoring them unless all inputs are missing.

By consistently favoring native PySpark SQL functions like **least** over custom Python UDFs, developers ensure that their code harnesses the full power of distributed computing. UDFs introduce significant overhead due to data serialization between the JVM and Python environments, which is entirely avoided when using optimized built-in functions. For large-scale data analysis, this performance difference can translate into hours of savings during complex ETL pipeline execution. Adhering to this principle is a foundational aspect of writing high-performance, production-ready PySpark code.

In conclusion, the **least** function is the authoritative tool for horizontal minimum aggregation. Its simplicity, combined with its profound performance benefits, solidifies its place as the standard solution for determining the minimum value across a specified set of numerical columns within a PySpark DataFrame. Mastering this function is a necessary step for any data engineer or scientist working with distributed data processing systems.