

Learning PySpark: Finding the Minimum Value by Group in a DataFrame

Authored by
Mohammed loot

November 11, 2025

RECOMMENDED CITATION

Mohammed loot (2025). *Learning PySpark: Finding the Minimum Value by Group in a DataFrame*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=16523>

Introduction to Grouped Minimum Calculation in PySpark

Analyzing massive datasets requires sophisticated techniques to derive meaningful summary insights. One of the most fundamental operations in big data processing is the calculation of summary statistics--such as the minimum, maximum, or average--across specific subgroups within the data. Working within the highly efficient [PySpark](#) framework, finding the minimum value for a numeric attribute, partitioned by the distinct categories defined in other columns, is a routine but critical task. This process falls under the general category of [Data Aggregation](#). [PySpark](#), as the Python API for [Apache Spark](#), offers exceptionally optimized functions that execute these calculations in a distributed and highly parallel manner, establishing it as the industry standard for handling petabytes of data across clustered environments. Mastering the correct application of grouping and minimum calculation functions is essential for any data engineer or scientist working on complex tasks, from financial risk modeling and anomaly detection to performance benchmarking and robust quality control systems.

The central data structure utilized within the [PySpark](#) ecosystem is the [DataFrame](#). Conceptually, this structure resembles a traditional SQL table, featuring named columns and distributed rows; however, it is specifically optimized for distributed, fault-tolerant computation across a cluster of machines. To effectively compute the minimum value based on specified groups, we rely on two pivotal methods inherent to the [DataFrame](#) object. First, the [groupBy](#) transformation is invoked, which logically organizes the data by partitioning it based on the values in the selected categorical column(s). Second, the powerful `agg()` method is applied, which takes the resulting partitioned blocks and executes the desired aggregation function--in this case, the calculation of the minimum value--on the specified target column.

This detailed guide will meticulously walk you through the two primary methodologies for achieving minimum value calculation by group. We will begin with the simplest case: grouping by a single categorical column, and then we will extend this logic to handle scenarios requiring simultaneous grouping across multiple columns. Throughout these examples, we will exclusively utilize the highly optimized built-in functions provided by the [pyspark.sql.functions](#) module, specifically focusing on the efficient [F.min\(\)](#) function. This approach ensures that the resulting code is not only concise and readable but also achieves peak performance when scaled across large clusters, which is vital for real-world big data operations.

To provide a quick reference, the following methods outline the high-level syntax required to calculate the minimum value by group within a [PySpark DataFrame](#):

Method 1: Calculate Minimum Grouped by One Column

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

```
#calculate minimum of 'points' grouped by 'team'  
df.groupBy('team').agg(F.min('points')).show()
```

Method 2: Calculate Minimum Grouped by Multiple Columns

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

```
#calculate minimum of 'points' grouped by 'team' and 'position'  
df.groupBy('team', 'position').agg(F.min('points')).show()
```

Setting Up the PySpark Environment and Sample Data

Before we can dive into the implementation of our [Data Aggregation](#) logic, it is absolutely essential to correctly initialize a functional Spark session and define the source dataset that will undergo analysis. For the purpose of maintaining clarity and ensuring these examples are easily reproducible, we will construct a small, representative dataset containing simulated basketball player statistics. This dataset is structured to include key categorical fields, namely the player's **team** and their **position**, alongside the numerical metric **points**, which is the column we intend to aggregate and find the minimum value from. This simple, yet effective, structure is perfectly suited to clearly demonstrate how PySpark's grouping logic correctly isolates the minimum scores within the specific, defined categories.

The initial setup phase requires importing the necessary library components, primarily focusing on the `SparkSession` constructor sourced from the `pyspark.sql` library. Once the session is successfully initialized--which establishes the connection to the underlying [Apache Spark](#) cluster--we then define our raw data structure. This data is typically represented as a list of lists or tuples in Python, and we must explicitly specify the corresponding column names (schema). The most crucial technical step is transforming this raw, inert Python data into a structured, distributed [DataFrame](#) using the `spark.createDataFrame()` method. While Spark often possesses the capability to infer the schema automatically, explicitly ensuring the input structure and data types are correct is paramount for avoiding errors and ensuring successful downstream analytical operations.

The following comprehensive code block details the entire setup process, illustrating the creation of the foundational [PySpark DataFrame](#). This DataFrame contains realistic information about various basketball players, structured into the three crucial columns: 'team', 'position', and the numerical metric 'points'. After the DataFrame is created, we invoke the `df.show()` command, which triggers the execution and displays the first 20 rows of the resulting distributed dataset, confirming the data structure is ready for aggregation.

```
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|position|points|
+----+-----+-----+
| A| Guard| 11|
| A| Guard| 8|
| A| Forward| 22|
| A| Forward| 22|
| B| Guard| 14|
| B| Guard| 14|
| B| Guard| 13|
| B| Forward| 7|
| C| Guard| 8|
| C| Forward| 5|
+----+-----+-----+
```

Understanding the Core Aggregation Methods

The process of aggregation in [PySpark](#) is implemented through a carefully orchestrated sequence of functional transformations. The primary goal is to shift the level of analysis from individual records to summarizing statistics about groups of records. This is achieved by combining the power of the `groupBy()` transformation with the flexibility of the `agg()` action. The `groupBy()` function, which we link to the official Spark documentation for [groupBy](#) operations, does not immediately compute a result; instead, it returns a `GroupedData` object. This intermediate object logically partitions the underlying [DataFrame](#) based on the unique values found in the key columns specified.

Once the data is logically grouped, the `agg()` method is called upon this `GroupedData` object. This method is the engine that executes the actual statistical computation across each of the partitions. When calculating the minimum, we pass a specific aggregation function, such as [F.min\(\)](#), as an argument to `agg()`. It is essential to understand that these functions--imported via `pyspark.sql.functions`--are highly optimized, native Spark functions designed for speed and efficiency in a [distributed systems](#) context. They are significantly faster and more resource-efficient than performing the same operation using User Defined Functions (UDFs) or standard Python list operations, which would break the optimized processing chain.

For instance, in the context of finding minimum performance scores, the system logically isolates all players belonging to 'Team A', then calculates the minimum 'points' score within that isolated subset. It repeats this identical operation independently and in parallel for 'Team B', 'Team C', and any other distinct group. The result of the `agg()` method is a new, much smaller [DataFrame](#), which contains the unique grouping columns and the newly calculated minimum value column. This transformation effectively reduces large amounts of granular data into concise, actionable summary statistics. This foundational understanding of the `groupBy` and `agg` mechanism is crucial for mastering not just minimum calculation, but all forms of [Data Aggregation](#) in [PySpark](#).

Example 1: Calculating Minimum Grouped by a Single Column

The most straightforward aggregation scenario involves determining the minimum value based on partitioning the data by a single characteristic. In the context of our player statistics, our concrete objective here is to isolate the lowest score recorded in the **points** column corresponding to each unique team identified in the **team** column. This task is efficiently accomplished by chaining the powerful [groupBy](#) function, where we specify the column 'team' as the key, directly with the `agg()` method, which executes the minimum calculation.

The syntax employed in [PySpark](#) is intentionally designed to be highly readable and intuitive: we explicitly instruct the Spark engine to logically partition the entire [DataFrame](#) based on the values in the 'team' column, and subsequently, for every distinct group created, calculate the minimum

value found in the 'points' column. A necessary prerequisite is the initial import of `pyspark.sql.functions` conventionally aliased as `F`, which allows the use of the highly concise and efficient notation like [F.min\(\)](#). This chained function call structure is a hallmark of idiomatic [PySpark](#) coding, guaranteeing high performance due to its optimization for distributed computing across large clusters, even when dealing with petabytes of raw input data.

The final output of this operation is a highly concise summary table. This resulting table will contain the primary grouping column ('team') and the newly calculated aggregation column, automatically named by default as `min(points)`. This specific transformation effectively reduces the detailed, player-level transactional data into clear, actionable summary statistics. These statistics provide immediate, high-level insight into the minimum performance threshold achieved by players associated with each unique team, which is invaluable for initial data profiling and comparative analysis across groups.

The following code snippet demonstrates the exact syntax used to calculate the minimum value in the **points** column, grouped exclusively by the values found in the **team** column:

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

```
#calculate minimum of 'points' grouped by 'team'  
df.groupBy('team').agg(F.min('points')).show()
```

```
+----+-----+  
|team|min(points)|  
+----+-----+  
| A| 8|  
| B| 7|  
| C| 5|  
+----+-----+
```

Based on the resulting output, we can definitively identify the minimum performance scores achieved by players within each team:

The minimum points value recorded for players on **team A** is **8**.

The minimum points value recorded for players on **team B** is **7**.

The minimum points value recorded for players on **team C** is **5**.

Example 2: Calculating Minimum Grouped by Multiple Columns

In real-world business intelligence and advanced analytical scenarios, it is frequently necessary to perform [Data Aggregation](#) based on the intersection of several categorical variables, providing a

much higher resolution of insight. For instance, determining the true minimum performance level might necessitate segmenting the data not just by 'team', but by the combination of 'team' and 'position' simultaneously. This requirement mandates grouping by multiple columns, a task that [PySpark](#) handles elegantly and seamlessly by simply passing a sequence or list of column names to the [groupBy](#) function instead of a single name.

When we define the composite grouping key as `('team', 'position')`, the underlying [DataFrame](#) is logically partitioned into distinct groups corresponding to every unique combination of these two values. For example, 'Team A' and 'Guard' form one specific, isolated group, while 'Team A' and 'Forward' form another entirely separate group. The subsequent aggregation step then applies the [F.min\(\)](#) function to calculate the minimum 'points' score achieved strictly within the boundaries of each of these finely defined sub-groups.

This approach of multi-column grouping yields a significantly greater analytical depth compared to a simple single-column aggregation. Instead of merely obtaining one single minimum score for an entire team, we gain clear visibility into the minimum scores achieved by specific positional roles within that team (e.g., minimum score for Team A Guards vs. minimum score for Team A Forwards). This granular, dimensional detail is indispensable for tasks that demand deep breakdowns, such as pinpointing performance discrepancies between various operational segments, defining highly specific minimum performance thresholds for individual roles, or performing complex dimensional modeling.

The following syntax is used to calculate the minimum value in the **points** column, grouped by the combined values of the **team** and **position** columns:

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

```
#calculate minimum of 'points' grouped by 'team' and 'position'
df.groupBy('team', 'position').agg(F.min('points')).show()
```

```
+---+-----+-----+
|team|position|min(points)|
+---+-----+-----+
| A| Guard| 8|
| A| Forward| 22|
| B| Guard| 13|
| B| Forward| 7|
| C| Forward| 5|
| C| Guard| 8|
+---+-----+-----+
```

From the resulting output, we can observe the precise, detailed minimum scores for each unique team-position combination:

The minimum points value for **Guards on team A** is **8**.

The minimum points value for **Forwards on team A** is **22**.

The minimum points value for **Guards on team B** is **13**.

The minimum points value for **Forwards on team B** is **7**.

The minimum points value for **Guards on team C** is **8**.

The minimum points value for **Forwards on team C** is **5**.

Beyond Minimum: Expanding Your PySpark Aggregation Toolkit

Calculating the minimum value by group is an essential and foundational skill for executing robust data analysis using the [PySpark](#) framework. By effectively combining the powerful [groupBy](#) transformation with the specialized [F.min\(\)](#) aggregation function, data analysts and engineers can efficiently summarize and profile huge distributed datasets. This efficiency holds true regardless of the complexity of the grouping key, whether it involves a single column or a complex, composite combination of multiple features. These methods form the indisputable foundation for generating descriptive statistics and conducting rigorous data profiling within the broader context of the [Apache Spark](#) ecosystem.

For practitioners who are committed to expanding their analytical expertise in [PySpark](#), it is critical to recognize the immense versatility of the ``agg()`` function. This function is not limited to calculating minimums; it fully supports a vast array of other critical statistical operations. These include functions such as maximum (``F.max()``), average (``F.avg()``), sum (``F.sum()``), and count (``F.count()``), all of which operate with the same high performance and [Data Aggregation](#) efficiency. Furthermore, advanced use cases allow for the calculation of multiple summary statistics simultaneously within a single invocation of the ``agg()`` call. This powerful capability significantly reduces computational overhead and optimizes resource utilization compared to performing separate, sequential aggregations for each desired metric, which is a common performance pitfall for beginners.

To continue building upon the foundational knowledge of grouping and aggregation demonstrated throughout this guide, we encourage exploration of related tutorials. Understanding how to handle missing data within aggregated groups or how to rename the resulting aggregated columns dynamically are the next logical steps in mastering large-scale data manipulation in [Apache Spark](#).