

How to Round PySpark Column Values to 2 Decimal Places

Authored by
stats writer

February 3, 2026

RECOMMENDED CITATION

stats writer (2026). *How to Round PySpark Column Values to 2 Decimal Places*.
PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=129326>

PySpark: Round Column Values to 2 Decimal Places

1. Introduction to Precision and Data Integrity in PySpark

Handling numerical data effectively is a core requirement for any robust data processing framework. In the context of big data analysis using [PySpark](#), ensuring numerical precision--specifically managing the number of decimal places--is vital for maintaining data integrity and facilitating accurate reporting. Many real-world datasets, particularly those derived from sensor readings, financial records, or scientific measurements, often contain highly precise floating-point numbers that are unnecessarily detailed for presentation or subsequent calculations. Therefore, techniques for controlled numerical rounding become indispensable tools in the data scientist's arsenal.

The Apache Spark environment, through its Python API, [PySpark](#), provides specialized functions within the `pyspark.sql.functions` module designed specifically for column-wise data transformation. When the objective is to standardize numerical columns by limiting their precision, such as rounding all values to exactly two decimal places, the built-in `round` function offers an optimized and straightforward solution. This approach ensures computational efficiency, which is paramount when dealing with large-scale [DataFrame](#) operations.

By leveraging the `round` function, users can specify the exact level of precision required, thereby simplifying complex numerical outputs without sacrificing meaningful information. This not only improves the readability of the resulting [DataFrame](#) but also adheres to common requirements in financial modeling and statistical analysis where figures must often conform to specific display standards, such as currency formats which typically mandate two decimal places. We will explore the specific syntax and practical application of this function to achieve precise data manipulation within [PySpark](#).

2. Understanding the Need for Data Rounding

The necessity for rounding arises primarily from two domains: data presentation and algorithmic stability. In terms of presentation, raw floating-point numbers often display many digits after the decimal point due to internal computer representations or the nature of their measurement. While accurate, these lengthy numbers can clutter reports and dashboards, making interpretation difficult for end-users. Rounding these values to a manageable precision, such as two decimal places, significantly enhances clarity and user experience, especially when dealing with variables like percentages, averages, or standardized scores.

Furthermore, from a data processing perspective, controlling precision can sometimes be critical for downstream tasks. For instance, comparing two floating-point numbers for equality can be

notoriously complex in computing due to minuscule differences stemming from binary representation. By explicitly rounding values before comparison or aggregation, we introduce standardization, mitigating potential issues related to floating-point arithmetic errors. This practice ensures that calculations based on these columns yield consistent and expected results across distributed computing environments like PySpark.

It is important to differentiate between merely truncating digits and proper mathematical rounding. Truncation simply cuts off the digits beyond the specified point, while rounding adheres to standard rules (e.g., rounding up if the next digit is five or greater, and down otherwise). The PySpark round function implements standard rounding rules, providing mathematically sound transformations essential for accurate data analysis. This function ensures that the resulting numbers represent the closest possible value at the specified level of precision.

3. The PySpark `round` Function: Syntax and Parameters

The method utilized in PySpark for achieving precise rounding is the `round` function, which is available after importing it from the `pyspark.sql.functions` library. This function is designed to operate on columns within a DataFrame, returning a new column object where every value has been subjected to the rounding operation. Its strength lies in its simplicity and efficiency, especially when applied across millions or billions of records handled by the Spark engine.

The basic syntax of the `round` function requires two primary arguments: the target column and the scale (or precision) to which the value should be rounded. The scale parameter dictates the number of digits that should remain after the decimal point. For instance, specifying a scale of `2` instructs PySpark to round the values to two decimal places. If a scale of `0` were used, the values would be rounded to the nearest integer. The functionality allows for both positive and negative scale arguments, where a negative scale rounds to the specified number of places to the left of the decimal point (e.g., rounding to the nearest ten or hundred).

To incorporate this function into a DataFrame transformation, it is typically used in conjunction with the `withColumn` method. The `withColumn` method is a fundamental operation in PySpark used to add a new column or replace an existing one. By chaining `withColumn` with the `round` function, we explicitly define the transformation logic for the new column, maintaining the immutability of the original DataFrame structure.

4. Step-by-Step Implementation Guide

The process for rounding a column in PySpark is highly standardized and involves only a few lines of code. The initial requirement is to ensure the `round` function is imported from the necessary library. This step is crucial, as the function is not automatically available in the default DataFrame methods, requiring explicit access via `pyspark.sql.functions`. Following importation, the

application involves targeting the specific column and defining the precision level.

Below is the precise syntax required to implement this operation. This code snippet assumes the existence of a `DataFrame` named `df` that contains a numerical column named `points` which we intend to standardize to two decimal places. The result is stored in a new `DataFrame`, `df_new`, which includes the original data plus the newly calculated rounded column.

```
from pyspark.sql.functions import round
```

```
#create new column that rounds values in points column to 2 decimal places  
df_new = df.withColumn('points2', round(df.points, 2))
```

This particular example creates a new column named `points2`, which effectively holds the values from the original `points` column after each observation has been processed by the `round` function, specifically set to a precision of 2. This method is highly recommended over iterative approaches (like UDFs in Python) for standard operations like rounding, as built-in functions are optimized for the distributed nature of the Spark execution engine. Using native Spark SQL functions ensures the operation is executed efficiently across all nodes of the cluster.

5. Example: Round Column Values to 2 Decimal Places in PySpark

To fully illustrate the practical application of the `round` function, we must first establish a representative `DataFrame`. This example uses simulated data representing average points scored by basketball teams, where the points are initially recorded with high precision (many decimal places). The objective is to standardize these scores for reporting purposes.

The setup process involves initializing a `PySpark SparkSession`, which is the entry point to programming Spark with the `Dataset` and `DataFrame` API. We then define the schema and populate the `DataFrame` using a list of tuples, ensuring the 'points' column is correctly interpreted as a numerical type suitable for the rounding operation.

The following comprehensive code block demonstrates the necessary steps to define and display our initial data structure. Notice the varying precision in the original 'points' column, which sets the stage for the required transformation:

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

```
#define data
```

```
data = ,
```

```
,
```


We achieve this by utilizing the `withColumn` method and passing the output of `round(df.points, 2)` as the column definition. The choice of `2` as the second argument is critical, as it dictates the required precision level. This process automatically handles the underlying distributed computation, ensuring that the rounding logic is applied consistently and rapidly across the entire dataset, regardless of its size.

The following syntax implements the rounding operation, creating the new column `points2`, and then immediately displays the transformed `DataFrame` to verify the results. Notice how the original `points` column is preserved, allowing for easy comparison between the raw and standardized data points:

```
from pyspark.sql.functions import round
```

```
#create new column that rounds values in points column to 2 decimal places
df_new = df.withColumn('points2', round(df.points, 2))
```

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

```
+-----+-----+-----+
| team| points|points2|
+-----+-----+-----+
| Mavs|18.3494| 18.35|
| Nets|33.5541| 33.55|
| Lakers|12.6711| 12.67|
| Kings|15.6588| 15.66|
| Hawks|19.3215| 19.32|
| Wizards|24.0399| 24.04|
| Magic|28.6843| 28.68|
| Jazz|40.0001| 40.0|
|Thunder|24.2365| 24.24|
| Spurs|13.9446| 13.94|
+-----+-----+-----+
```

7. Analyzing the Results and Data Consistency

Upon viewing the resulting `DataFrame`, `df_new`, it is evident that the new column `points2` successfully incorporates the rounded values. The primary objective of standardizing the data to two decimal places has been achieved. A detailed analysis of the output confirms that the `PySpark round` function correctly implemented conventional mathematical rounding rules, ensuring data

consistency and accuracy.

For instance, examining the entry for 'Kings', the original value was 15.6588. Since the third decimal digit (8) is greater than or equal to 5, the second decimal digit (5) is rounded up, resulting in 15.66 in the **points2** column. Conversely, the 'Nets' score of 33.5541 rounds down to 33.55 because the third decimal digit (4) is less than 5. This adherence to strict mathematical rules is crucial for applications where precision tolerances are important.

The following list provides specific examples demonstrating how the rounding operation was executed for different types of fractional values:

The value **18.3494** was rounded up to **18.35**.

The value **33.5541** was rounded down to **33.55**.

The value **12.6711** was rounded down to **12.67**.

The value **24.0399** was rounded up to **24.04**.

The value **40.0001** was rounded down to **40.0** (note that PySpark often displays trailing zeros for decimal types, but in this output, it simplified to 40.0 as it is still interpreted as a precise float or decimal type internally).

8. Comparison with Alternative PySpark Rounding Methods

While the `round` function is the standard choice for general mathematical rounding, PySpark also provides alternative functions tailored for specific numerical requirements. One notable alternative is the `bround` function (Banker's Rounding), also imported from `pyspark.sql.functions`. Banker's Rounding differs from conventional rounding specifically when the digit to be rounded is exactly 5, in which case it rounds toward the nearest even number. This method is statistically preferred in some disciplines as it tends to reduce upward bias over a large series of calculations.

Furthermore, users might consider using casting operations if the primary goal is not rounding but simply enforcing a specific data type and scale. For instance, casting a high-precision float to a `Decimal(10, 2)` type will implicitly truncate or round the value depending on the Spark configuration. However, explicit use of the `round` or `bround` function is generally recommended when precise control over the rounding methodology is required, as casting often relies on implicit system rules which can vary.

In most standard business intelligence and data preparation pipelines, the simple `round` function provides the necessary precision control and transparency. It is the most intuitive function for ensuring that numerical columns conform to common display requirements, such as currency or standardized test scores, which demand two decimal places. Users should select the method that best aligns with the mathematical standards required by their specific application domain.

9. Conclusion: Ensuring Data Quality through Precision

The ability to efficiently manage and standardize numerical precision is a cornerstone of effective data transformation in big data environments. By employing the native [PySpark `round` function](#) within a `withColumn` operation, data engineers and analysts can quickly and reliably format data within large [DataFrames](#). This technique is computationally efficient and ensures adherence to precise mathematical [rounding](#) rules.

The methodology detailed here provides a high-performance solution for a common data manipulation task: forcing numerical values into a standardized format, such as two decimal places. This not only enhances the quality and readability of the data but also prepares it optimally for subsequent analytical modeling or reporting stages, where consistency is paramount. Mastering these basic, yet powerful, [PySpark SQL functions](#) is essential for achieving high levels of data integrity and trustworthiness in any Spark-based project.

For further exploration into the capabilities of the PySpark data manipulation functions, users are encouraged to consult the official documentation, which provides extensive details on all available mathematical and analytical transformations.

Further Resources on PySpark Data Manipulation

You can find the complete documentation for the [PySpark `round` function](#) on the Apache Spark official website.

The following tutorials explain how to perform other common tasks in [PySpark](#):