

# Convert String to Integer in PySpark (With Example)

Authored by  
**stats writer**

November 16, 2025

## RECOMMENDED CITATION

stats writer (2025). *Convert String to Integer in PySpark (With Example)*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=92504>

## Introduction to String to Integer Conversion in PySpark

PySpark, the Python API for Apache Spark, is an indispensable tool for large-scale data processing and analysis. When working with complex, unstructured, or semi-structured datasets, it is common to encounter numerical data stored inadvertently as a string type. Proper data typing is essential for accurate calculations, aggregations, and performance optimization within a DataFrame.

The core process for converting a column from a string representation to a numeric format, specifically an integer, involves leveraging the `cast()` function provided by the PySpark SQL module. This function allows developers to explicitly change the data types of columns within a PySpark DataFrame, ensuring compatibility for mathematical operations.

The general syntax required to execute this type conversion is straightforward and relies on importing the necessary type definition, such as IntegerType, from the `pyspark.sql.types` module. This method is robust and preferred for maintaining data integrity during transformation pipelines in big data environments. The following snippet illustrates the fundamental syntax utilized for this conversion:

```
from pyspark.sql.types import IntegerType
```

```
df = df.withColumn('my_integer', df.cast(IntegerType()))
```

In the exemplary code above, a new column named **my\_integer** is generated. This column inherits its values from the original **my\_string** column, but those values are explicitly converted using the `cast()` operation applied in conjunction with the IntegerType definition. This approach ensures that the original string column remains unmodified, allowing for safer, non-destructive transformations. The following sections demonstrate this syntax through a practical, step-by-step example.

## Understanding PySpark Data Types and the cast() Method

Effective data manipulation in PySpark necessitates a deep understanding of how data types are handled within the Spark ecosystem. Unlike standard Python, Spark relies on its own set of internal SQL types optimized for distributed computing. When loading data, especially from external sources like CSV or JSON files, columns that contain numeric information are often mistakenly inferred or explicitly loaded as a string type, which restricts analytical capabilities and can lead to significant performance bottlenecks.

The primary mechanism for changing a column's data type is the `cast()` function. This is a transformation applied directly to a PySpark column object, instructing Spark to attempt conversion to a specified target type. If the conversion is successful (e.g., converting the string '123' to the

integer 123), the operation proceeds smoothly across all partitions. However, if the conversion fails (e.g., trying to cast the non-numeric string 'N/A' to an integer), Spark gracefully handles the error by inserting a **null** value into the output column for that specific record, maintaining the integrity of the transformation flow.

When defining the target type, it is crucial to use the specific type objects imported from ``pyspark.sql.types``. For integer conversion, we use `IntegerType`. This object specifies the 32-bit signed integer format expected by Spark's SQL engine. The ``withColumn`` function is used to integrate this transformation, allowing us to either overwrite the existing column or, as recommended, create a new column based on the casted values of the source column, thereby providing an auditable history of the data transformation.

## Prerequisites: Setting up the PySpark Environment and DataFrame

To demonstrate the practical application of the string-to-integer conversion, we must first establish a Spark session and create a representative sample `DataFrame`. Our example simulates a common scenario in sports analytics where raw score data, representing points scored by various basketball teams, has been loaded as string values due to preliminary data ingestion processes. This setup ensures we start with the exact condition requiring type conversion.

The process begins with the necessary imports and the initialization of the Spark session. We then define our sample data, consisting of 'team' identifiers and corresponding 'points'. It is important to note that the values in the 'points' column are deliberately defined as strings within the data list, forcing the resulting `DataFrame` schema to recognize them as non-numeric.

This code block initializes the environment, loads the data into a PySpark `DataFrame` named `df`, and displays the resulting dataset for visual confirmation:

```
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|points|
+----+-----+
| A| 11|
| B| 19|
| C| 22|
| D| 25|
| E| 12|
| F| 41|
| G| 32|
| H| 20|
+----+-----+
```

The resulting `df` is now ready for inspection. The next step is to programmatically verify the schema to confirm that the 'points' column, despite containing only numerical digits, is currently represented as a string, preventing any direct mathematical calculations.

## Step 1: Inspecting Initial Data Types (The Necessity of Conversion)

To formally confirm the current structure and data types assigned by Spark, we utilize the built-in `dtypes` attribute of the DataFrame. This attribute returns a list of tuples, where each tuple contains the column name and its corresponding Spark SQL type, providing a clear map of the DataFrame's structure.

Executing `df.dtypes` is a crucial diagnostic step. If we attempted aggregation functions like `sum()` or `avg()` on a string column containing numbers, the operation would either fail outright or treat the values as categorical data, yielding incorrect results. It is important to distinguish between the appearance of data (e.g., '11') and its underlying logical type (e.g., string).

The output below clearly demonstrates the current schema, highlighting the specific column requiring conversion:

## #check data type of each column df.dtypes

As evidenced by the output, the **points** column is currently stored as a `string` type. This confirmation establishes the clear and immediate need to convert this column to an integer format before any meaningful quantitative analysis can be performed. The next step addresses this requirement directly using the appropriate PySpark transformation.

## Step 2: Implementing the String-to-Integer Conversion using IntegerType()

The actual transformation is executed using the `withColumn` method in PySpark, which facilitates column manipulation. We aim to create a new column, named `points_integer`, which will hold the converted values derived from the original `points` column. This approach is highly recommended for data lineage tracking.

The syntax involves three key components: the `withColumn` function, the selection of the source column, and the application of the `cast()` method paired with `IntegerType`. We first ensure that `IntegerType` is imported from the `pyspark.sql.types` module, providing the necessary type object for the cast operation.

Upon execution, Spark processes the `cast` instruction across all records. If all string values are valid numerical representations, the transformation is seamless. The resulting DataFrame, shown below, now contains both the original string column and the new, numerically functional integer column:

```
from pyspark.sql.types import IntegerType
```

```
#create integer column from string column  
df = df.withColumn('points_integer', df.cast(IntegerType()))
```

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

```
+----+-----+-----+  
|team|points|points_integer|  
+----+-----+-----+  
| A| 11| 11|  
| B| 19| 19|  
| C| 22| 22|  
| D| 25| 25|
```

```
| E| 12| 12|  
| F| 41| 41|  
| G| 32| 32|  
| H| 20| 20|  
+---+-----+
```

While the visual output confirms that the numbers appear identical in both columns, the underlying storage and processing capabilities of the new `points_integer` column have been fundamentally altered, paving the way for optimized numerical computation. The final step is to formally verify this internal schema change.

### Step 3: Verification of the Updated DataFrame Schema

Following any significant data transformation, particularly a type conversion, it is absolutely essential to verify the resulting schema. We once again employ the `df.dtypes` command to inspect the updated structure of the DataFrame, ensuring that the new column possesses the intended data types.

Successful conversion is confirmed by observing the type associated with the new column, which must be `int`. This step is crucial for maintaining confidence in the data pipeline, especially in production environments where downstream systems rely on correctly typed data for aggregation and model training.

Executing the schema check yields the following output:

```
#check data type of each column  
df.dtypes
```

The output unequivocally confirms that the `points_integer` column has been correctly assigned the `int` type. This validates that we have successfully achieved the goal of creating an integer column from the source `string` column using the PySpark `cast()` operation with `IntegerType`. The data is now ready for complex numerical processing within the Spark ecosystem.

### Conclusion and Best Practices for Data Type Management

Mastering data type conversion is a cornerstone of effective big data processing in `PySpark`. Converting columns from string to integer using the `cast(IntegerType())` method ensures that numerical data is handled optimally by the Spark engine, leading to improved performance, reduced memory footprint, and enhanced accuracy in statistical analysis and machine learning.

pipelines.

When working with production-scale datasets, it is vital to adhere to best practices regarding data typing. Always profile your data types initially using `df.dtypes` or `df.printSchema()` to identify conversion needs proactively. Furthermore, when performing conversions, consider the appropriate numerical type size. While `IntegerType` (32-bit signed integer) is often sufficient, data types like `ShortType` (16-bit) or `LongType` (64-bit) may be necessary depending on the scale and magnitude of the numerical values being processed, optimizing resource consumption.

Finally, always be prepared to handle potential conversion errors resulting from non-conforming data. The standard `cast()` handles exceptions gracefully by introducing nulls, but for robust data pipelines, users might employ techniques like filtering out corrupted records or utilizing SQL functions such as `regexp_replace` to clean non-numeric characters before attempting the final cast. By adhering to these guidelines, data practitioners can ensure smooth and reliable data transformation workflows.