

PySpark: Check Data Type of Columns in DataFrame

Authored by
Mohammed loot

November 11, 2025

RECOMMENDED CITATION

Mohammed loot (2025). *PySpark: Check Data Type of Columns in DataFrame*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=16511>

Why Data Type Inspection is Crucial in PySpark

The ability to inspect and verify the schema of a [DataFrame](#) is fundamental when performing data engineering tasks using [PySpark](#). Unlike traditional Python objects where types are sometimes inferred dynamically, Spark relies heavily on explicitly defined or correctly inferred [data types](#) for optimized processing across a distributed cluster. If the types are incorrectly handled--for example, treating a numerical column as a string--performance degradation or unexpected data loss during transformations can occur.

In a big data environment, knowing the precise **data type** assigned to each column is essential for debugging and ensuring efficiency. PySpark provides straightforward methods to retrieve this schema information using the built-in **dtypes** attribute of the DataFrame object. These methods allow developers to quickly assess the structure of their data, whether they need to check a single column or the entire dataset.

We will explore two primary methods available in PySpark for checking column data types, followed by practical examples demonstrating their implementation and interpreting the resulting outputs.

Setting Up the Environment and Sample DataFrame

Before demonstrating the type-checking methods, we must first establish a working environment and create a sample [DataFrame](#). This foundational step utilizes the [SparkSession](#) object, which serves as the entry point for all PySpark functionality. The following Python code snippet demonstrates the creation of a simple dataset containing information about sports teams, including numerical and categorical fields.

This example dataset is designed to showcase how PySpark automatically infers different types, such as **string** for textual data and **bigint** for integer data, which is standard behavior when creating DataFrames from basic Python lists without an explicit schema definition.

```
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|conference|points|assists|
+----+-----+-----+-----+
| A| East| 11| 4|
| A| null| 8| 9|
| A| East| 10| 3|
| B| West| null| 12|
| B| West| null| 4|
| C| East| 5| 2|
+----+-----+-----+-----+
```

Method 1: Inspecting the Data Type of a Single Column

When performing targeted column operations, such as filtering based on a specific field or casting it to a new type, it is highly efficient to verify only that column's data type. This approach avoids unnecessary processing overhead associated with retrieving the schema for the entire [DataFrame](#).

To check the type of a single column, we utilize the **df.dtypes** attribute, but wrap it within the Python **dict()** function. This conversion allows us to treat the schema output as a key-value mapping, where the column name acts as the key, providing immediate access to its associated type.

We can use the following syntax to check the data type of the **conference** column in the DataFrame:

```
#return data type of 'conference' column
dict(df.dtypes)

'string'
```

The resulting output, **'string'**, confirms that the **conference** column is treated as a StringType by Spark SQL. This is expected since the column contains categorical text data.

To check the data type of another specific column, such as **points**, simply replace **conference** with the desired column name while maintaining the dictionary access structure:

```
#return data type of 'points' column  
dict(df.dtypes)
```

```
'bigint'
```

The output confirms that the **points** column has a **bigint** data type, which is PySpark's standard inference for integers when the values exceed the range of a standard 32-bit integer, or simply based on default configuration settings.

Understanding the `df.dtypes` Output Structure

The core component used in both methods is the DataFrame attribute [.dtypes](#). It is important to understand what this attribute returns before it is converted into a dictionary. When called directly, **df.dtypes** returns a standard Python list of tuples.

Each tuple within this list contains two elements: the name of the column (a string) and its associated Spark SQL data type (also a string). For instance, an excerpt of the raw output might look like **('team', 'string')**. While iterating through this list is possible, converting it to a Python dictionary is generally preferred for rapid, index-based access, especially when querying the type of a known column name.

The dictionary conversion method provides a clean, Pythonic way to handle schema lookup, significantly improving code readability and performance compared to looping through the list of tuples manually. This technique showcases how PySpark interoperates seamlessly with standard Python data structures to simplify complex data manipulation tasks.

Method 2: Checking Data Types Across the Entire DataFrame

When initiating work on a new dataset, or after performing several complex transformations (such as joins or aggregations), it is often necessary to review the schema in its entirety. This comprehensive approach ensures that all columns align with expected types before proceeding with data aggregation or persistence.

The simplest and most direct method for checking all column data types is to call the **.dtypes** attribute on the DataFrame object without any modification or conversion. This returns the complete list of tuples representing the schema.

We can use the following syntax to check the data type of all columns in the DataFrame:

```
#return data type of all columns  
df.dtypes
```

The result clearly maps every column name to its inferred or defined type. This output is critical for validating the initial data load and for confirming that transformations have not inadvertently changed the schema in an undesirable way.

From this output, we can deduce the following:

The **team** column has a data type of [string](#).

The **conference** column has a data type of **string**.

The **points** column has a data type of [bigint](#).

The **assists** column has a data type of **bigint**.

Implications of PySpark Data Types (BigInt vs. Integer)

It is vital to understand the difference between standard Python types and the specialized SQL types used by PySpark. When PySpark infers the schema from raw Python data, it often defaults to the safest, most comprehensive type. For instance, the numerical columns, **points** and **assists**, were assigned the [bigint](#) type, which corresponds to **LongType** in Spark's internal structure.

A **bigint** utilizes 8 bytes (64 bits) of storage, allowing it to hold very large integer values. While this offers maximum flexibility, if you know your data will never exceed the range of a standard 4-byte integer (`IntType`), explicitly casting the column to **integer** can result in significant memory savings across a large-scale cluster. Checking the [data type](#) before persisting or processing the data allows developers to optimize storage and computational resources.

Similarly, text-based columns are typically inferred as **string** (`StringType`), which is the standard representation for variable-length character sequences in Spark SQL. If optimization is necessary, developers might consider alternative types, but **string** is usually the correct choice for categorical text data like team names or conference locations.

Conclusion and Additional Resources

Mastering the inspection of **data types** is a fundamental skill in PySpark development. Whether you require a quick check of a single column using the dictionary method or a full schema review using the raw **df.dtypes** attribute, these techniques ensure data integrity and facilitate necessary type casting for optimized performance. Always verify your schema after data ingestion or complex transformation steps to maintain control over your distributed dataset.

The following tutorials explain how to perform other common tasks in PySpark, building upon this core understanding of DataFrame schema management: