

# PySpark: Check Data Type of Columns in DataFrame

Authored by  
**stats writer**

November 16, 2025

## RECOMMENDED CITATION

stats writer (2025). *PySpark: Check Data Type of Columns in DataFrame*.

PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=92542>

## The Importance of Schema Inspection in PySpark

When working with large datasets in the Apache Spark environment, particularly using the PySpark API, understanding the underlying structure of your data is paramount. This structure, known as the schema, dictates how Spark interprets, processes, and allocates memory for each value within a DataFrame. Incorrect data type assignment can lead to processing errors, inefficient resource utilization, or unexpected results during complex transformations and aggregations. Therefore, regularly verifying the data type of columns is a fundamental step in any robust data pipeline.

PySpark offers straightforward and efficient methods for inspecting the data type of columns within a DataFrame. Whether you need a quick check on a single column or a comprehensive overview of the entire dataset structure, the API provides the necessary tools. We will explore two primary methods here: targeted inspection for a single column and a holistic approach for checking all columns simultaneously. These techniques ensure that the data types--such as **string**, **integer**, **bigint**, or **double**--align perfectly with your analytical requirements.

Effective schema management is especially critical when dealing with data ingested from external sources, where Spark might infer the schema based on limited samples, sometimes leading to non-optimal type assignments. By actively inspecting and confirming these types, data engineers can preemptively resolve type mismatch issues. The following sections detail the precise syntax and implementation for querying column data types, starting with the methods themselves before moving into practical, runnable examples.

### Core Concepts: Understanding Spark DataFrames and Data Types

A Spark DataFrame conceptually resembles a table in a relational database or a data frame in R/Python, but with optimizations tailored for distributed computing. Each column in this structure is assigned a specific data type, which dictates the kind of values it can hold and the operations that can be performed on it. PySpark utilizes the concept of `dtypes`--a property of the DataFrame object--to hold this schema information as a list of tuples, where each tuple contains the column name and its corresponding type (e.g., `('column_name', 'type')`).

Understanding Spark's internal data type nomenclature is vital. For instance, while Python might use standard `int`, Spark often uses **IntegerType**, **LongType** (represented as 'bigint' in the output), or **StringType**. When checking the data types in PySpark, the output will reflect these internal Spark SQL types. The efficiency of operations like joins, filters, and aggregations is heavily dependent on having numerical data stored in appropriate numerical types (like **bigint** or **double**) rather than as strings.

The techniques demonstrated below leverage the `df.dtypes` attribute, which serves as the central access point for schema information. While `df.printSchema()` offers a visually formatted output

of the [schema](#), directly accessing `df.dtypes` provides a structured Python list that is far more useful for programmatic analysis and extraction, allowing users to manipulate or query the type information directly within their code.

## Method 1: Inspecting a Single Column's Data Type

To efficiently retrieve the [data type](#) for only one specific column, we can treat the `df.dtypes` list of tuples as a key-value structure. Although `df.dtypes` is inherently a list, casting it to a Python **dictionary** allows for quick lookups using the column name as the key. This approach is highly useful when debugging an issue related to a single field or confirming the type conversion of a newly created column.

The core of this method involves two steps: first, retrieving the list of all column types using `df.dtypes`, and second, converting this list into a standard Python `dict`. Once converted, standard dictionary indexing syntax (e.g., `dict[column]`) can be used to instantly return the corresponding Spark SQL data type as a string.

Here is the concise syntax for this operation. This method focuses solely on the column specified, making it efficient for large [DataFrames](#) where inspecting the entire schema is unnecessary overhead.

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

## Implementation Details: Retrieving the Data Type Dictionary

The code snippet above, `dict(df.dtypes)`, performs a crucial transformation. The `df.dtypes` property returns a list of tuples, like `[(column, type)]`. Python's built-in `dict()` constructor can accept a list of two-item tuples and automatically transform it into a dictionary where the first element of the tuple becomes the key and the second becomes the value. In this case, the column name becomes the key, and the Spark data type string becomes the value.

Once this conversion is complete, accessing on the resulting dictionary performs a rapid  $O(1)$  lookup. This is significantly faster and cleaner than manually iterating through the list of tuples to find the matching column name, especially when dealing with [DataFrames](#) containing hundreds of columns. This programmatic accessibility is a key advantage for integrating schema checks into automated validation scripts.

While this method is simple, users must ensure the column name used as the key exactly matches the case and spelling of the column name in the DataFrame, as key lookup in Python dictionaries is case-sensitive. This method will fail gracefully (raise a `KeyError`) if the specified column does

not exist, alerting the user to potential issues in column naming or availability.

## Method 2: Comprehensive Schema Inspection Using `.dtypes`

When the goal is to review the structure of the entire DataFrame--perhaps immediately after data ingestion or schema modification--the simplest and most direct approach is to access the `.dtypes` attribute without any further manipulation. This method returns the raw list of tuples containing all column names and their corresponding data type strings, providing a complete overview of the schema.

This technique is generally preferred when initially exploring a dataset or when migrating code, as it quickly reveals whether any columns were unexpectedly inferred as **string** instead of a numerical or date type. Unlike `printSchema()`, which provides a formatted, print-only output, `df.dtypes` returns a Python object (the list of tuples) that can be stored in a variable, processed further, or compared against expected schema definitions.

The output is a canonical representation of the PySpark schema. Below is the syntax, demonstrating its elegant simplicity. This method is the foundation upon which more complex schema validation frameworks are built in large-scale PySpark applications.

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

## Practical Setup: Creating the Sample DataFrame

To illustrate these two methods effectively, we must first establish a reproducible PySpark DataFrame. The following setup code initializes a **SparkSession**, defines the sample data, specifies the desired column names, and creates the distributed structure. This example dataset includes a mix of string and numerical data, simulating typical operational data where teams, conferences, points, and assists are tracked.

Notice that the numerical columns ('points' and 'assists') contain some `None` values in the Python data definition. When PySpark infers the schema from this data, it defaults to the most appropriate Spark SQL type based on the presence of integers. Since PySpark 3.x, integer types are often inferred as **LongType** (which appears as 'bigint' in the output) to accommodate potential scaling, unless explicitly specified otherwise in the schema definition.

The code block below generates the required environment and displays the resulting DataFrame using `df.show()`, confirming the structure before we proceed with schema interrogation.

```
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|
```

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

### Example 1: Targeted Inspection of the 'conference' Column

In this first practical application, we use Method 1 to specifically isolate and verify the data type assigned to the **conference** column. This column holds categorical text data, and we anticipate that PySpark correctly inferred it as a **string** type. Using the dictionary conversion approach ensures that only the required type information is retrieved.

By executing the following code snippet against our sample DataFrame `df`, we confirm that the mechanism correctly extracts the type associated with the specified column name. This confirms that even if the source data had mixed types, Spark's inference mechanism or an explicitly defined schema dictated the final type of **string** for this field.

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

```
'string'
```

The output, `'string'`, confirms our expectation that the **conference** column is handled as a standard text field by Spark SQL. This result is crucial; if this column had mistakenly been inferred as an integer or some other numerical type (which might happen in poorly structured CSV files), subsequent SQL operations or filtering based on string values would fail or produce runtime errors.

## Expanding the Scope: Checking Other Specific Columns

The elegance of Method 1 lies in its reusability. To check the data type of any other column, one simply replaces the key in the dictionary lookup. Let us now apply this technique to a numerical column, **points**, which contains integer values and nulls, to observe how Spark manages numerical data types.

We must remember that PySpark's default inference often maps standard Python integers to **LongType** (represented as `bigint`), which is a 64-bit signed integer. This is done to prevent overflow issues when dealing with large datasets typical of the Spark environment. If one requires 32-bit integers, an explicit schema must typically be provided during DataFrame creation.

Executing the lookup for the **points** column reveals its inferred type:

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

```
'bigint'
```

The resulting type, `'bigint'` (or **bigint**), confirms that the **points** column is ready for standard mathematical operations while also having sufficient capacity to store large integer values, reinforcing the reliability of Spark's default schema inference for numerical data.

## Example 2: Reviewing the Full DataFrame Schema

Finally, we demonstrate Method 2, which provides a comprehensive view of the entire DataFrame structure. This is often the starting point for any data quality or validation task in PySpark. By calling `df.dtypes` directly, we receive a list that fully describes the schema, allowing for immediate comparison against expected standards.

The output is an ordered list of tuples, reflecting the column order in the DataFrame. Each tuple

explicitly pairs the column name with its determined Spark SQL data type. This holistic view is necessary when performing transformations that affect multiple columns simultaneously or when ensuring cross-compatibility between different datasets.

Here is the execution and the resulting full schema output:

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

Analyzing the output confirms the schema derived by PySpark during DataFrame creation. We can precisely determine the type of every column:

The **team** column has a data type of **string**, suitable for categorical labels.

The **conference** column has a data type of **string**, also for categorical grouping.

The **points** column has a data type of **bigint** (LongType), appropriate for large integer counts.

The **assists** column also has a data type of **bigint**, mirroring the numerical type of points.

This complete schema inspection is vital for the validation phase of any data processing workflow, providing the necessary assurance that the data structure is sound before proceeding to complex analytical queries.

## Conclusion: Ensuring Data Integrity in PySpark

Mastering the art of schema inspection in PySpark is a critical skill for any data professional utilizing the framework. By employing the two straightforward methods detailed--the targeted dictionary lookup for specific columns and the direct `df.dtypes` access for the full schema--you gain immediate control and insight into how Spark is interpreting your data.

These techniques are invaluable for diagnostics, performance tuning, and ensuring that downstream transformations and machine learning models receive data structured precisely as expected. Consistent verification of data types prevents silent errors, maintains data quality, and ensures the efficient execution that Apache Spark is renowned for.

Incorporating these schema checks into standard operating procedures guarantees robust and reliable data pipelines, allowing engineers to focus on analysis rather than dealing with unexpected type coercions or casting failures.