

# How to Create an Empty PySpark DataFrame with Defined Columns

Authored by  
**stats writer**

February 7, 2026

## RECOMMENDED CITATION

stats writer (2026). *How to Create an Empty PySpark DataFrame with Defined Columns*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=129728>

Creating an empty `DataFrame` in `PySpark` with specific column names is a fundamental task for data professionals initiating large-scale data pipelines or transformations. This process allows developers to establish a robust, predefined structure--or schema--before any data arrives. This structure is essential because `PySpark` is highly reliant on knowing the data types and column names in advance for optimized execution across a distributed cluster.

The standard method involves leveraging the `SparkSession` and the powerful `createDataFrame()` method. While this method typically accepts data (like an RDD or a list of rows) and a schema, we can intentionally pass an empty data source (an empty list or an empty RDD) while providing a fully defined, custom schema. By explicitly defining the column names and their associated data types using classes like `StructType` and `StructField`, we achieve complete control over the resultant `DataFrame` structure immediately upon creation.

This approach is invaluable for development, testing, and situations where data loading is deferred until later stages of an application. It ensures that subsequent operations expecting specific columns and types will execute smoothly, regardless of whether the initial dataset is populated. The subsequent sections will detail the required imports, the critical role of schema definition, and the exact syntax used to execute this procedure efficiently in a `PySpark` environment.

## PySpark: Create Empty DataFrame with Column Names

### Setting Up the PySpark Environment

Before we can successfully generate an empty `DataFrame`, it is necessary to initialize the Spark environment and import the specific libraries required for defining the data structure. The core component for interaction is the `SparkSession`, which serves as the entry point to programming Spark with the `Dataset` and `DataFrame` API. We also require components from the `pyspark.sql.types` module to define the complex structure of our schema.

The first step involves importing the `SparkSession` class itself. If a session already exists (common in environments like Databricks or Jupyter Notebooks), `builder.getOrCreate()` ensures we utilize the existing session, preventing redundant initialization. If no session is active, a new one is instantiated, handling all necessary configurations behind the scenes.

Furthermore, defining a `DataFrame` requires explicit declaration of its columns and their corresponding data types. For this purpose, we specifically import `StructType` and `StructField`, which are the building blocks of any complex schema in `PySpark`. This foundational setup guarantees that the subsequent code for creating the empty `DataFrame` runs smoothly, adhering to the distributed computing principles of Apache Spark.

## Defining the Schema: StructType and StructField

The critical difference between creating a populated `DataFrame` and an empty one lies in the mandatory, explicit definition of the `schema`. When Spark reads data, it usually tries to infer the schema, but since we are providing no data, we must tell Spark exactly what to expect. This is achieved using two primary classes from `pyspark.sql.types`.

The `StructType` class represents the entire row structure of the `DataFrame`. It acts as a container for a list of `StructField` objects. Think of `StructType` as the blueprint for the table itself. Conversely, each `StructField` object defines a single column. It requires three key pieces of information: the column name (a string), the data type (e.g., `StringType()`, `IntegerType()`), and a boolean indicating whether the column allows null values (nullability).

By assembling a list of `StructField` definitions and wrapping them within a `StructType` object, we create a comprehensive schema object. This schema is then passed directly to the `createDataFrame()` function, instructing `PySpark` on how to label and interpret the columns, even though the data payload itself is nonexistent. This precise definition ensures type safety and prevents unexpected errors during subsequent transformations.

## Core Syntax for Empty DataFrame Creation

The following code block demonstrates the concise and standard syntax used to create an empty `DataFrame` named `df`. Notice how we combine the setup of the `SparkSession` with the definition of the schema and the final creation call. While the original example included the creation of an empty RDD, modern `PySpark` practice often simplifies this by passing an empty list as the data source parameter.

The key line here is the definition of `my_columns`, where we specify the structure: a column named 'team' (string), 'position' (string), and 'points' (float). Each column is set to be nullable (`True`), which is typically the default behavior but is explicitly stated here for clarity and robust definition.

You can use the following syntax to create an empty `PySpark DataFrame` with specific column names:

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

```
from pyspark.sql.types import StructType, StructField, StringType, FloatType
```

```
#create empty RDD (Optional, but shown for compatibility)
empty_rdd=spark.sparkContext.emptyRDD()
```

```
#specify colum names and types
my_columns=

#create DataFrame with specific column names
df=spark.createDataFrame(, schema=StructType(my_columns))
```

This particular example successfully initializes a DataFrame called **df**. Crucially, it establishes three columns: **team** (intended for categorical data), **position** (also textual), and **points** (designed to hold numerical, floating-point values). The structure is now fixed, awaiting the ingestion of data from various sources.

## Detailed Implementation Example and Visualization

To demonstrate the functionality, the following example expands upon the core syntax by immediately visualizing the resulting structure. When an empty DataFrame is displayed using the `df.show()` command, the output confirms the presence of the defined column headers (team, position, points) but contains no rows, visually confirming that the structure exists independently of the data.

The ability to define an empty, structured DataFrame is highly valuable in scenarios such as iterative processing or pipeline initialization. For instance, if you are reading data in chunks or conditionally loading data, having an empty DataFrame ready allows you to append results or fill the structure seamlessly without worrying about schema mismatch errors, which are common when combining or merging data sources in distributed computing.

The following example shows how to use this syntax in practice, confirming that the columns are properly established:

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

from pyspark.sql.types import StructType, StructField, StringType, FloatType

#create empty RDD (This step is technically optional if passing to createDataFrame)
empty_rdd=spark.sparkContext.emptyRDD()

#specify colum names and types
my_columns=

#create DataFrame with specific column names
df=spark.createDataFrame(, schema=StructType(my_columns))
```

```
#view DataFrame
df.show()

+----+-----+-----+
|team|position|points|
+----+-----+-----+
+----+-----+-----+
```

As evidenced by the output of `df.show()`, an empty PySpark DataFrame has been successfully created. While no records are present, the column headers: **team**, **position**, and **points**, are clearly displayed, confirming the correct application of the defined schema.

## Verifying the DataFrame Schema with `printSchema()`

While `df.show()` confirms the column names, the critical step in validation is ensuring that the data types were correctly applied according to the definitions provided by the StructType. The `df.printSchema()` method provides a hierarchical, detailed view of the schema, which is essential for debugging and ensuring type compatibility, especially when integrating with other systems or databases.

This method outputs the structure starting from the root (the DataFrame itself), listing each field (column), its assigned data type (e.g., string, float), and its nullability status (whether it can contain missing values). For our example, this validation confirms that 'points' is indeed defined as a float, allowing for decimal numeric data, while 'team' and 'position' are defined as strings, suitable for textual identifiers.

We can also use the following syntax to view the schema of the DataFrame, which confirms the data types and nullability:

```
#view schema of DataFrame
df.printSchema()
```

```
root
|-- team: string (nullable = true)
|-- position: string (nullable = true)
|-- points: float (nullable = true)
```

From this comprehensive output, we can deduce the following structural properties:

The **team** field is designated as a string type, confirming it is designed to hold text data.

The **position** field is also defined as a string type, suitable for textual labels or categories.

The **points** field is successfully defined as a float type, ensuring it can handle numerical values that may include decimal points.

## Understanding Data Types and Nullability

The choice of data type in the `StructField` definition is crucial, as PySpark relies heavily on these types for efficient memory management and distributed processing. Using the wrong type can lead to performance degradation or casting errors when data is finally loaded. Common types include `StringType`, `IntegerType`, `FloatType`, `DoubleType`, `BooleanType`, and complex types like `ArrayType` or `MapType`.

The third parameter in the `StructField` constructor, the boolean value (e.g., `True`), dictates the column's nullability. Setting it to `True` means the column can accept missing values (nulls), while setting it to `False` enforces non-null constraints. Although PySpark often does not strictly enforce non-null constraints until specific write operations, explicitly defining nullability is a best practice for clean schema management and data governance.

It is important to consult the official Apache Spark documentation to see the complete list of available data types, as selecting the most appropriate type ensures maximum performance and compatibility across the Spark ecosystem. For instance, using `IntegerType` instead of `StringType` for numerical IDs saves significant resources during large-scale operations.

## Conclusion: Leveraging Structured Empty DataFrames

Creating an empty, schema-defined DataFrame is not merely an academic exercise; it is a critical preparatory step in robust data engineering using PySpark. This technique allows for the early definition of data contracts, ensuring consistency across disparate data sources and providing a reliable target structure for subsequent data ingestion, transformation, and analysis tasks.

By mastering the use of `StructType` and passing an empty list to the `createDataFrame()` method, developers gain the flexibility to build complex data pipelines where the structure is finalized before the data volume necessitates distributed computation. This practice significantly improves code maintainability and reduces the likelihood of schema evolution errors in production environments.

**Note:** You can find a complete list of data types that you can specify for columns in a PySpark DataFrame by referring to the official `pyspark.sql.types` documentation.

## Further PySpark Exploration

Having established the foundation for creating structured DataFrames, you can now explore various other common tasks essential for data processing in a distributed environment.

Understanding the manipulation, joining, and aggregation of these DataFrames is the next logical step. These tutorials explain how to perform other common tasks in PySpark:

How to efficiently join multiple DataFrames using different join types.

Methods for filtering and selecting data based on complex conditions.

Techniques for using Window functions for advanced analytical operations.

ARABPSYCHOLOGY.COM