

# Read CSV File into PySpark DataFrame (3 Examples)

Authored by  
**stats writer**

November 16, 2025

## RECOMMENDED CITATION

stats writer (2025). *Read CSV File into PySpark DataFrame (3 Examples)*.

PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=92539>

Working with large datasets often requires importing external data sources, and the **Comma-Separated Values (CSV)** format remains one of the most ubiquitous standards for data exchange. When leveraging the scalability and distributed computing power of PySpark, the primary tool for loading this tabular data into memory is the `spark.read.csv()` function. This powerful function allows users to ingest complex, multi-gigabyte files efficiently into a structured DataFrame, preparing the data for large-scale analysis and transformation. Understanding the core parameters of this method is essential for any data engineer or data scientist utilizing the Apache Spark ecosystem.

The process of reading a CSV file into a PySpark DataFrame can be customized based on the structure of the source file, such as the presence of a header row or the type of delimiter used. We will explore three fundamental approaches to demonstrate the flexibility of the reader API. These examples showcase how to handle basic ingestion, incorporate header information, and manage non-standard delimiters, providing a comprehensive foundation for handling diverse CSV inputs.

### Three Essential Methods for Reading CSV Data

The `spark.read.csv()` function provides several optional arguments that significantly influence how the data is parsed and structured upon loading. The three methods outlined below represent the most frequently encountered scenarios when dealing with raw CSV data ingestion. Mastering these techniques ensures that the imported data accurately reflects the source structure, preventing downstream errors during data processing.

#### Method 1: Basic CSV File Reading (Default Settings)

```
df = spark.read.csv('data.csv')
```

#### Method 2: Reading CSV File with Header Enforcement

```
df = spark.read.csv('data.csv', header=True)
```

#### Method 3: Reading CSV File with Specific Delimiter Definition

```
df = spark.read.csv('data.csv', header=True, sep=';')
```

The subsequent sections delve into detailed, practical applications of each of these methods, illustrating the setup required--including initializing a SparkSession--and analyzing the structure of the resulting DataFrame after ingestion.

## Example 1: Basic CSV Reading Without Header Definition

In the first scenario, we assume a basic CSV file where we initially omit the explicit specification of header information. This demonstrates the default parsing behavior of PySpark. We will utilize a hypothetical file named **data.csv**, which contains standard team statistics, including team designation, total points, and assists count. It is crucial to observe how PySpark handles the first row of data when the header option is not set.

Consider the contents of our sample file, **data.csv**, structured using standard comma delimiters:

**team, points, assists**

**'A', 78, 12**

**'B', 85, 20**

**'C', 93, 23**

**'D', 90, 8**

**'E', 91, 14**

To read this file using the minimal configuration, we initialize the Spark environment and invoke the `spark.read.csv()` function, passing only the file path. This approach relies entirely on PySpark's default assumptions regarding data structure, which often results in the first row being treated as data rather than metadata, and automatic column naming conventions being applied. This is generally suitable only when the source file truly lacks descriptive headers.

The following syntax illustrates the initial setup and the subsequent data loading:

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

#read CSV into PySpark DataFrame using default settings
df = spark.read.csv('data.csv')

#view resulting DataFrame structure and data
df.show()
```

```
+----+-----+-----+
|_c0|_c1|_c2|
+----+-----+-----+
|team| points| assists|
| 'A'| 78| 12|
| 'B'| 85| 20|
| 'C'| 93| 23|
```

```
| 'D'| 90| 8|  
| 'E'| 91| 14|  
+----+-----+-----+
```

## Understanding PySpark's Default Ingestion Behavior

When the **header** option is omitted or explicitly set to **False**, PySpark makes a critical assumption: that the input file does not contain a header row. Consequently, the first line of the CSV file is interpreted as the first record of data, leading to the column names being automatically generated. These generated names follow a sequential pattern, typically starting with **\_c0** for the first column, **\_c1** for the second, and so on. This mechanism ensures that every column has a unique identifier, even if meaningful descriptive names are absent.

In the output shown above, the actual headers (team, points, assists) are present in the first data row under the generic column names **\_c0**, **\_c1**, and **\_c2**. While this default behavior is predictable, it often necessitates an additional step of cleaning or renaming columns if the user intended for the first row to define the schema. Furthermore, all columns are initially loaded as generic strings, requiring explicit type casting later if numerical operations are needed, unless the **inferSchema** option is used (which is typically slower but can determine data types automatically).

The use of these default settings is generally discouraged for production datasets that possess clear metadata headers. Instead, it is preferable to utilize the available options within the `spark.read.csv()` function to accurately define the file structure during the loading phase itself, optimizing both readability and initial data types.

### Example 2: Reading CSV File with Explicit Header Recognition

The most common requirement for data ingestion is to recognize and utilize the first row of the file as the column headers. PySpark facilitates this by using the **header=True** parameter. This setting instructs the reader to skip the first line when reading data records but to use its values when defining the schema for the resulting DataFrame. This method significantly improves the usability and clarity of the resulting structure, aligning the DataFrame schema with the expectations set by the source file.

We will reuse the identical **data.csv** file from Example 1 to demonstrate the impact of this critical parameter:

```
team, points, assists  
'A', 78, 12  
'B', 85, 20
```

'C', 93, 23

'D', 90, 8

'E', 91, 14

The following syntax is used to initialize the `SparkSession` and load the CSV, specifically setting the header parameter to **True**. This simple addition transforms how the data is interpreted, making the subsequent data manipulation cleaner and more intuitive, as we can now reference columns by their meaningful names (e.g., 'team' instead of '\_c0').

Implementation demonstrating header enforcement:

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

```
#read CSV into PySpark DataFrame, recognizing the header row
df = spark.read.csv('data.csv', header=True)
```

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

```
+----+-----+-----+
|team| points| assists|
+----+-----+-----+
| 'A'| 78| 12|
| 'B'| 85| 20|
| 'C'| 93| 23|
| 'D'| 90| 8|
| 'E'| 91| 14|
+----+-----+-----+
```

As evident from the output, setting **header=True** ensures that PySpark correctly maps the column names from the first row of the CSV file to the DataFrame schema. This is the recommended practice for reading structured data, as it preserves semantic meaning and simplifies all subsequent operations carried out using the DataFrame API.

### Example 3: Reading CSV File with Specific Delimiter Specification

While the term "CSV" implies comma separation, many systems use alternative characters such as semicolons, tabs, or pipes to delimit fields, particularly in regions where the comma is commonly used as a decimal separator. PySpark is fully equipped to handle these variations through the **sep**

(separator) argument within the `spark.read.csv()` function. This flexibility is crucial for handling international data formats or outputs generated by proprietary systems.

For this example, imagine that our `data.csv` file now uses semicolons (;) instead of commas to separate the values:

**team; points; assists**

'A'; 78; 12

'B'; 85; 20

'C'; 93; 23

'D'; 90; 8

'E'; 91; 14

If we were to attempt reading this file using the default comma delimiter, the entire row would be mistakenly parsed as a single column. To correctly parse this structure, we must explicitly define the separator using the `sep` parameter, setting it to the semicolon character (;). We also maintain `header=True` to ensure the column names are correctly assigned.

The required syntax for handling the semicolon delimiter:

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

```
#read CSV into PySpark DataFrame, specifying header and semicolon delimiter
df = spark.read.csv('data.csv', header=True, sep=';')
```

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

```
+----+-----+-----+
|team| points| assists|
+----+-----+-----+
| 'A'| 78| 12|
| 'B'| 85| 20|
| 'C'| 93| 23|
| 'D'| 90| 8|
| 'E'| 91| 14|
+----+-----+-----+
```

By successfully utilizing the `sep` argument, `PySpark` correctly recognized the semicolon as the field separator, leading to a properly structured `DataFrame` where the data is segmented into the

appropriate columns (team, points, assists). This capability highlights the robustness of the `spark.read.csv()` function in adapting to various input file standards.

## Advanced Considerations for Robust CSV Loading

While the examples above cover the most frequent use cases, handling enterprise-level or dirty data often requires utilizing additional options provided by the CSV reader. Parameters such as **inferSchema**, **nullValue**, **dateFormat**, and **quote** are vital for creating a truly resilient loading process. The **inferSchema=True** option, for instance, instructs Spark to perform an extra pass over the data to deduce the correct data types (e.g., Integer, Double, Date) instead of treating everything as strings, although this adds overhead to the loading time.

For complex scenarios involving malformed records, the **mode** parameter is crucial. By default, Spark uses **PERMISSIVE** mode, which inserts nulls for fields that cannot be parsed. However, options like **FAILFAST** (to abort upon encountering bad data) or **DROPMALFORMED** (to silently drop corrupted rows) offer granular control over error handling. A truly professional ingestion pipeline often requires predefining a schema using **StructType** and **StructField** objects to bypass schema inference entirely, providing maximal speed and deterministic data types.

Ultimately, the strength of the PySpark CSV reader lies in its wide array of configuration options, allowing developers to tailor the data ingestion to the specific constraints and quality of the source CSV file. By combining the basic methods demonstrated here with these advanced configurations, users can ensure clean, fast, and reliable loading of massive datasets into the Spark ecosystem for subsequent processing and analysis.