# Spark by {Examples} (https://sparkbyexamples.com/)

Spark # (http://sparkbyexamples.com/)

PySpark   (https://sparkbyexamples.com/pyspark-tutorial/)

Hive   (https://sparkbyexamples.com/apache-hive-tutorial/)

HBase   (https://sparkbyexamples.com/apache-hbase-tutorial/)

Kafka   (https://sparkbyexamples.com/apache-kafka-tutorials-with-examples/)

FAQ's   (https://sparkbyexamples.com/spark-interview-questions/)

More ⌄   (https://sparkbyexamples.com/) 🔍

# PySpark – Create DataFrame with Examples

👤 NNK (https://sparkbyexamples.com/author/admin/) -

📁 PySpark (https://sparkbyexamples.com/category/pyspark/)

You can manually create a **PySpark DataFrame** using `toDF()` and `createDataFrame()` methods, both these function takes different signatures in order to create DataFrame from existing RDD, list, and DataFrame.

You can also create PySpark DataFrame from data sources like TXT, CSV, JSON, ORV, Avro, Parquet, XML formats by reading from HDFS, S3, DBFS, Azure Blob file systems e.t.c.

**Related:**

- Fetch More Than 20 Rows & Column Full Value in DataFrame (https://sparkbyexamples.com/spark/spark-fetch-more-than-20-rows-full-column-value/)
- Get Current Number of Partitions of Spark DataFrame (https://sparkbyexamples.com/spark/spark-get-current-number-of-partitions-of-dataframe/)

Finally, PySpark DataFrame also can be
created by reading data from RDBMS
Databases and NoSQL databases.

In this article, you will learn creating
DataFrame by some of these methods
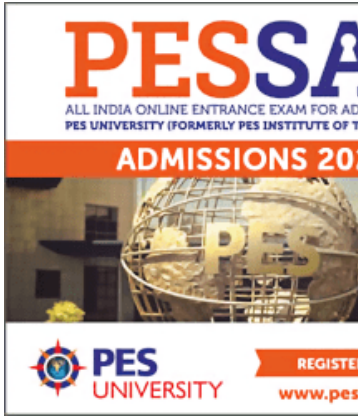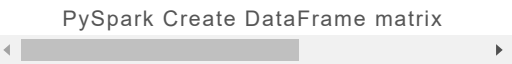with PySpark examples.



## Table of Contents

- Create DataFrame from RDD
  - toDF()
  - createDataFrame()
- Create DataFrame from the list of
  data
- Create DataFrame from Data sources
  - Creating from CSV file
  - Creating from TXT file
  - Creating from JSON file
- Other sources (Avro, Parquet, ORC
  e.t.c)

| SPARKSESSION | F |
|---|---|
| createDataFrame(rdd) | t |
| createDataFrame(dataList) | t |
| createDataFrame(rowData,columns) | |
| createDataFrame(dataList,schema) | |

PySpark Create DataFrame matrix

In order to create a DataFrame from a
list we need the data hence, first, let's
create the data and the columns that
are needed.

```
columns = ["language","users_co
data = [("Java", "20000"), ("Py
```

# 1. Create DataFrame from RDD

One easy way to manually create PySpark DataFrame is from an existing RDD. first, let's [create a Spark RDD (https://sparkbyexamples.com/spark/different-ways-to-create-spark-rdd/)](https://sparkbyexamples.com/spark/different-ways-to-create-spark-rdd/) from a collection List by calling [parallelize() (https://sparkbyexamples.com/pyspark/pyspark-parallelize-create-rdd/)](https://sparkbyexamples.com/pyspark/pyspark-parallelize-create-rdd/) function from [SparkContext (https://sparkbyexamples.com/spark/spark-sparkcontext/)](https://sparkbyexamples.com/spark/spark-sparkcontext/). We would need this **rdd** object for all our examples below.

```
spark = SparkSession.builder.ap
rdd = spark.sparkContext.parall
```

## 1.1 Using toDF() function

PySpark RDD's toDF() method is used to create a DataFrame from existing RDD. Since RDD doesn't have columns, the DataFrame is created with default column names "_1" and "_2" as we have two columns.

```
dfFromRDD1 = rdd.toDF()
dfFromRDD1.printSchema()
```

printschema() yields the below output.

```
root
 |-- _1: string (nullable = tru
 |-- _2: string (nullable = tru
```

If you wanted to provide column names to the DataFrame use `toDF()` method with column names as arguments as shown below.

```
columns = ["language","users_co
dfFromRDD1 = rdd.toDF(columns)
dfFromRDD1.printSchema()
```

This yields schema of the DataFrame
with column names.

```
root
 |-- language: string (nullable
 |-- users: string (nullable =
```

By default, the datatype of these
columns infers to the type of data. We
can change this behavior by supplying
schema
(https://sparkbyexamples.com/pyspark/
pyspark-structtype-and-structfield/),
where we can specify a column name,
data type, and nullable for each
field/column.

## 1.2 Using createDataFrame() from SparkSession

Using createDataFrame() from
SparkSession
(https://sparkbyexamples.com/tag/spark
session) is another way to create
manually and it takes rdd object as an
argument. and chain with toDF() to
specify name to the columns.

```
dfFromRDD2 = spark.createDataFr
```

## 2. Create DataFrame from List Collection

In this section, we will see how to create PySpark DataFrame from a list. These examples would be similar to  t we have seen in the above section  RDD, but we use the list data  ct instead of "rdd" object to create aFrame.

### Using createDataFrame() from SparkSession

ing `createDataFrame()` from  rkSession  is another way to  ate PySpark DataFrame manually, it  s a list object as an argument. and  n with `toDF()` to specify names to columns.

```
fFromData2 = spark.createDataF
```

### Using createDataFrame() with the  w type

ateDataFrame() has another  ature in PySpark which takes the  ction of Row type and schema for  mn names as arguments. To use  first we need to convert our "data"  ct from the list to list of Row.

```
rowData = map(lambda x: Row(*x)
dfFromData3 = spark.createDataF
```

## 2.3 Create DataFrame with schema

If you wanted to specify the column names along with their data types, you should create the StructType schema first and then assign this while creating a DataFrame.

```python
from pyspark.sql.types import S
data2 = [("James","","Smith","3
    ("Michael","Rose","","40288
    ("Robert","","Williams","42
    ("Maria","Anne","Jones","39
    ("Jen","Mary","Brown","","F
  ]

schema = StructType([ \
    StructField("firstname",Str
    StructField("middlename",St
    StructField("lastname",Stri
    StructField("id", StringTyp
    StructField("gender", Strin
    StructField("salary", Integ
  ])

df = spark.createDataFrame(data
df.printSchema()
df.show(truncate=False)
```

This yields below output.

```
root
 |-- firstname: string (nullable = true)
 |-- middlename: string (nullable = true)
 |-- lastname: string (nullable = true)
 |-- id: string (nullable = true)
 |-- gender: string (nullable = true)
 |-- salary: integer (nullable = true)

+---------+----------+--------+-----+------+------+
|firstname|middlename|lastname|id   |gender|salary|
+---------+----------+--------+-----+------+------+
|James    |          |Smith   |36636|M     |3000  |
|Michael  |Rose      |        |40288|M     |4000  |
|Robert   |          |Williams|42114|M     |4000  |
|Maria    |Anne      |Jones   |39192|F     |4000  |
|Jen      |Mary      |Brown   |     |F     |-1    |
+---------+----------+--------+-----+------+------+
```

## 3. Create DataFrame from Data sources

In real-time mostly you create DataFrame from data source files like CSV, Text, JSON, XML e.t.c.

PySpark by default supports many data formats out of the box without importing any libraries and to create DataFrame you need to use the appropriate method available in DataFrameReader class.

## 3.1 Creating DataFrame from CSV

Use `csv()` method of the `DataFrameReader` object to create a DataFrame from CSV file. you can also provide options like what delimiter to use, whether you have quoted data, date formats, infer schema, and many more. Please refer PySpark Read CSV into DataFrame (https://sparkbyexamples.com/pyspark/pyspark-read-csv-file-into-dataframe/)

```
df2 = spark.read.csv("/src/reso
```

## 3.2. Creating from text (TXT) file

Similarly you can also create a DataFrame by reading a from Text file, use `text()` method of the DataFrameReader to do so.

```
df2 = spark.read.text("/src/res
```

## 3.3. Creating from JSON file

PySpark is also used to process semi-structured data files like JSON format. you can use `json()` method of the DataFrameReader to read JSON file into DataFrame. Below is a simple example.

```
df2 = spark.read.json("/src/res
```

Similarly, we can create DataFrame in PySpark from most of the relational databases which I've not covered here and I will leave this to you to explore.

# 4. Other sources (Avro, Parquet, ORC, Kafka)

We can also create DataFrame by reading Avro, Parquet, ORC, Binary files and accessing Hive and HBase table, and also reading data from Kafka which I've explained in the below articles, I would recommend reading these when you have time.

- PySpark Read Parquet file into DataFrame (https://sparkbyexamples.com/pyspark/pyspark-read-and-write-parquet-file/)
- DataFrame from Avro source (https://sparkbyexamples.com/spark/using-avro-data-files-from-spark-sql-2-4/)
- DataFrame by Streaming data from Kafka (https://sparkbyexamples.com/spark/spark-streaming-kafka-consumer-example-in-json-format/)

The complete code can be downloaded from GitHub (https://github.com/spark-examples/spark-scala-examples/blob/master/src/main/scala/com/sparkbyexamples/spark/dataframe/CreateDataFrame.scala)

Happy Learning !!

---

**Share this:**

TAGS: **DATAFRAME
(HTTPS://SPARKBYEXAMPLES.COM/TAG/DATAFRA
ME/)**, **SPARK.CREATEDATAFRAME
(HTTPS://SPARKBYEXAMPLES.COM/TAG/CREATED
ATAFRAME/)**, **SPARKCONTEXT
(HTTPS://SPARKBYEXAMPLES.COM/TAG/SPARKCO
NTEXT/)**, **SPARKCONTEXT.PARALLELIZE
(HTTPS://SPARKBYEXAMPLES.COM/TAG/SPARKCO
NTEXT-PARALLELIZE/)**, **SPARKSESSION
(HTTPS://SPARKBYEXAMPLES.COM/TAG/SPARKSE
SSION/)**, **TODF()
(HTTPS://SPARKBYEXAMPLES.COM/TAG/TODF/)**

## NNK

## (Https://Sparkbyexamples.Com/Author/Admin/)

(https://sp
arkbyexa
mples.co
m/author/
admin/)

SparkByExamples.com is a Big Data and Spark
examples community page, all examples are simple and
easy to understand and well tested in our development
environment Read more ..
(https://sparkbyexamples.com/about-sparkbyexamples/)

❯ THIS POST HAS 3 COMMENTS

### Anonymous

**4 FEB 2021**      **REPLY**

* indicates: its passing
list as an argument

### Anonymous

**24 NOV 2020**      **REPLY**

regular expression for
arbitrary column names

### sagar      **30 JUL 2020**      **REPLY**

What is significance of * in below
dfFromData2 =
spark.createDataFrame(data).toDF(*col
umns)

## Leave a Reply

Enter your comment here...

ark-check-string-column-has-numeric-
values/)

Spark Check Column Data Type is
Integer or String
(https://sparkbyexamples.com/spark/sp
ark-check-column-data-type-is-integer-
or-string/)

ark-check-string-column-has-numeric-
values/)

Spark Check Column Data Type is
Integer or String
(https://sparkbyexamples.com/spark/sp
ark-check-column-data-type-is-integer-
or-string/)