

## Q) What is PySpark :-

PySpark is an interface for Apache Spark in Python. It is a library.

## Q) Setting Up :-

• Create a virtual environment.

• Install pyspark in it :-

[pip install pyspark]

It is also install py4j automatically.

## Q) SparkSession :-

• When we want to work with PySpark, we always need to start a Spark Session.

To create a Spark Session, we need to code as follows :-

```
[from pyspark.sql import SparkSession  
spark = SparkSession.builder.appName ("Prac...  
...tise").getOrCreate ()]
```

Name of the Spark session that we create

This function actually creates the session, and returns the session object.

We save the session in a variable (spark here)

• We can run this code multiple times, and it creates a session (takes a bit of time when 1st time executed).

⇒ The Spark Session gets created in memory/RAM.  
When running in local, there will be 1 cluster.  
On cloud, there will be Master & multiple clusters.

## ⇒ Reading data with PySpark :-

SparkSession object ; file name/path that is to be read.

[df\_spark = spark.read.csv("filename.extension")]

storing the data ; here, we have many options like  
that is being read ; csv, json, .. Depending on the  
file to be read, that function should  
be used appropriately .

⇒ This method will take all the data as raw data, and will store them as column data (even the column headings will be considered data), and will by itself give arbitrary naming to column as headings.  
E:- c0, -c1.

⇒ To read by taking the 1st row as row header names, we need to write as follows:-

[df\_spark = spark.read.option('header', 'true').csv("filename")]

Due to this option() method, we can set the 1st row as header row.

⇒ 1st parameter : Key of option, (which is 'header' here)  
2nd param: Value of the key.

⇒ Displaying the data being read :-

[ df\_spark.show() ]

⇒ show() method displays the data stored in the PySpark dataframe object, in tabular format.

⇒ [type (df\_spark) → pyspark.sql.DataFrame.DataFrame]

⇒ This shows what type of object, in which the read data is being stored.

⇒ Thus, we see, it is a PySpark SQL DataFrame

⇒ head() method :-

Takes one integer parameter (n) :-

[ df\_spark.head(n) ]

⇒ n > 1 :- Returns a list of first (n) rows.

⇒ n = 1 :- Returns a single row object.

⇒ printSchema() method :-

⇒ This method actually prints the schema of the

