# Ingeniería en sistemas computaciones

Datos masivos



**Profesor**

Jose Christian Romero

# Unit 1

Exam 1

**Alumnos:**

Marquez Millan Seashell Vanessa

Galaviz Lona Oscar Eduardo

Desarrollo

**1-.SPARK session**

Only start SPARK with the comand

```
SPARK-SHELL
```



**2-. File Netflix Stock CSV**

First need import the library, and is important you have the dataFrame in this address "/home/"name computer"/"name dataFrame" because the comand stearchh the archive here, then only print the data types

```scala
import org.apache.spark.sql.SparkSession

val spark = SparkSession.builder().getOrCreate()

val df = spark.read.option("header",
"true").option("inferSchema","true")csv("Netflix_2011_2016.csv")
df
df.printSchema()
```

```
scala> df.printSchema()
root
 |-- Date: timestamp (nullable = true)
 |-- Open: double (nullable = true)
 |-- High: double (nullable = true)
 |-- Low: double (nullable = true)
 |-- Close: double (nullable = true)
 |-- Volume: integer (nullable = true)
 |-- Adj Close: double (nullable = true)


scala> df
res47: org.apache.spark.sql.DataFrame = [Date: timestamp, Open: double ... 5 more fields]

scala> 
```
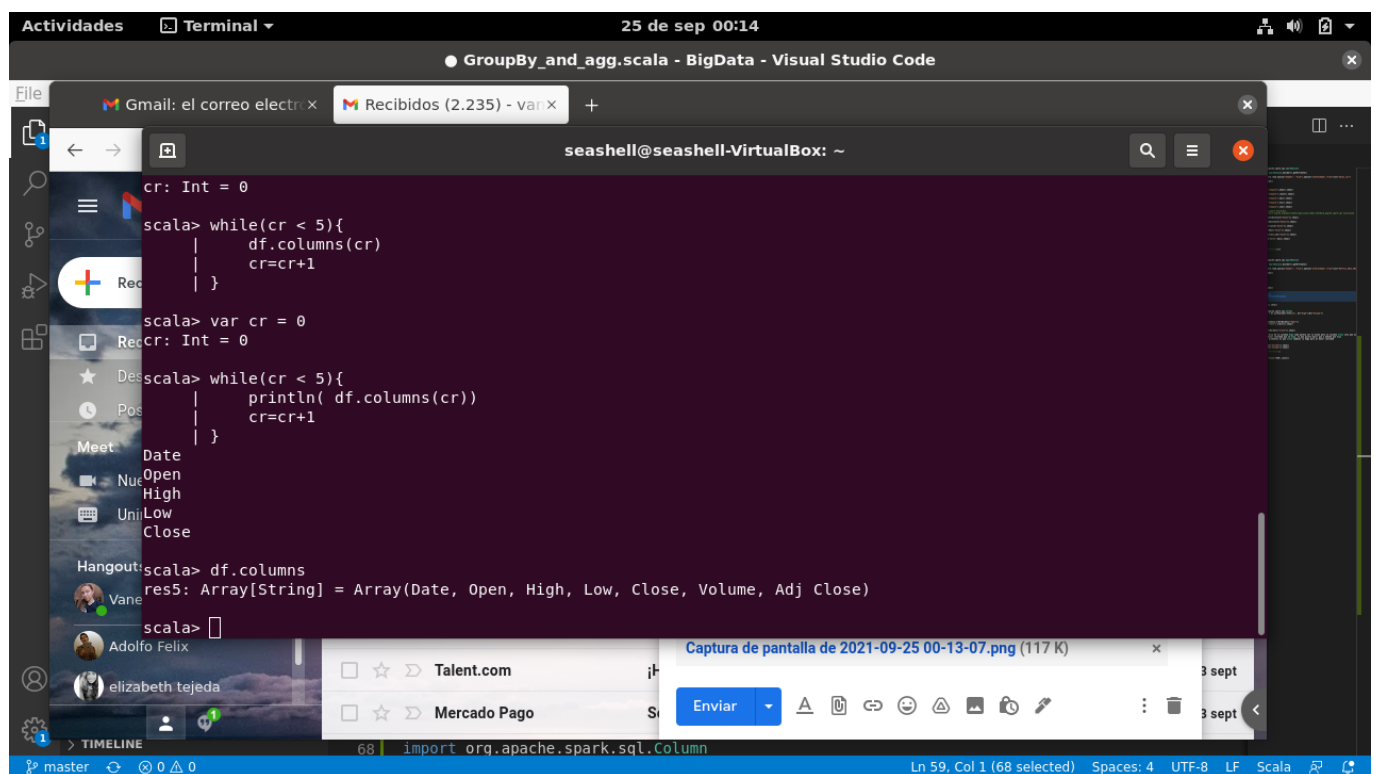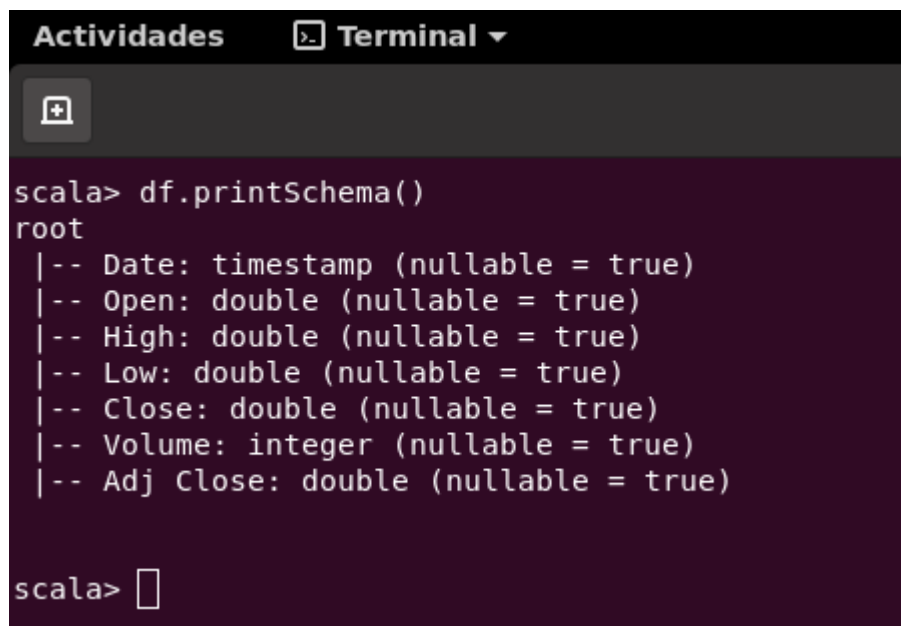
### 3-.Names columns

Only start SPARK with the comand

```
df.columns
```



### 4-.Schema

For know thw schema only need the dataFrame and the next reserverd word, is for can you know the structure
and the types of each column

```
df.printSchema()
```

```
Actividades      Terminal ▾

  ⊞                                           4 / 9

scala> df.printSchema()
root
 |-- Date: timestamp (nullable = true)
 |-- Open: double (nullable = true)
 |-- High: double (nullable = true)
 |-- Low: double (nullable = true)
 |-- Close: double (nullable = true)
 |-- Volume: integer (nullable = true)
 |-- Adj Close: double (nullable = true)


scala> []
```

### 5-Print first 5 columns

For that need create a variable inthis case cr and i say is equals to 0, then use while for create a bucle ever cr is less than 5, print the columns and increase cr if i don't do that is a infinite bucle

```
var cr = 0
while(cr < 5){
    println( df.columns(cr))
    cr=cr+1
}
```

## 6-.Uses describe ()

That comand is for knows more informations about the dataFrame, statistical data

```
df.describe().show()
```

**7-.Create new DataFrame with new column**

We need create a new dataFrame for can to do some modification so here to make a new column with the relationship of column High and Volume

```scala
val newData = df.withColumn("HVRatio", df("High")/df("Volume"))
newData.show()
```



**8-.Max Open**

We need know the date of the maxium data, so first we order the column Open and save in maxp then select Date of de maxp but only the first row

```scala
val maxp = newData.orderBy(desc("Open"))
maxp.select("Date").limit(1).show()
```

```
5|
|2015-07-09 00:00:00|        664.300011|        670.919975|        659.999992|        670.089996|16076900|        95.727142|4.173192437596800..
.|
|2015-07-01 00:00:00|663.6400219999999|666.6699980000001|652.5300219999999|        655.449982|14699300|        93.635712|4.535386025184873E-
5|
|2015-06-18 00:00:00|662.5599980000001|667.3900150000001|        660.579979|        663.200012| 8962800|        94.742859|7.446222330075424E-
5|
|2015-06-16 00:00:00|        659.700012|        669.249977|        655.840004|        666.910004|16043300|        95.272858|4.171523171666676..
.|
|2015-07-02 00:00:00|        657.990013|        659.389992|        652.500008|        658.31002|11053000|        94.044289|5.965710594408758E-
5|
|2015-07-06 00:00:00|        654.309982|664.5000150000001|653.3799740000001|661.9999849999999|11808300|        94.571426|5.627397804933818E-
5|
|2015-07-08 00:00:00|        654.299995|657.9699860000001|        645.990005|        654.549988|12990600|        93.507141|5.064969947500500..
.|
+------------------+----------------+----------------+----------------+----------------+--------+----------------+------------------
-+
only showing top 20 rows

maxp: Unit = ()

scala>

scala> val maxp = newData.orderBy(desc("Open"))
maxp: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [Date: timestamp, Open: double ... 6 more fields]

scala> maxp.select("Date").limit(1).show()
+------------------+
|              Date|
+------------------+
|2015-07-14 00:00:00|
+------------------+

scala>
```

### 9-.Meaning Close in DataFrame

```
newData.orderBy(desc("Close")).show()
```

### 10-.Maximum and minimum of Volume

This is only to know thw first row the most big and the most lowest, and oly select the volume and your minium or maximun

```
df.select(max("Volume")).show()
df.select(min("Volume")).show()
```

```
Actividades      Terminal ▾                         24 de sep 22:56                                          ⬇ ◀)) 🔋 ▾

⊞                                    seashell@seashell-VirtualBox: ~                                    🔍  ≡  ✕

  at org.apache.spark.sql.catalyst.analysis.checkAnalysis$class.checkAnalysis(checkAnalysis.scala:85)
  at org.apache.spark.sql.catalyst.analysis.Analyzer.checkAnalysis(Analyzer.scala:95)
  at org.apache.spark.sql.catalyst.analysis.Analyzer$$anonfun$executeAndCheck$1.apply(Analyzer.scala:108)
  at org.apache.spark.sql.catalyst.analysis.Analyzer$$anonfun$executeAndCheck$1.apply(Analyzer.scala:105)
  at org.apache.spark.sql.catalyst.plans.logical.AnalysisHelper$.markInAnalyzer(AnalysisHelper.scala:201)
  at org.apache.spark.sql.catalyst.analysis.Analyzer.executeAndCheck(Analyzer.scala:105)
  at org.apache.spark.sql.execution.QueryExecution.analyzed$lzycompute(QueryExecution.scala:57)
  at org.apache.spark.sql.execution.QueryExecution.analyzed(QueryExecution.scala:55)
  at org.apache.spark.sql.execution.QueryExecution.assertAnalyzed(QueryExecution.scala:47)
  at org.apache.spark.sql.Dataset$.ofRows(Dataset.scala:78)
  at org.apache.spark.sql.Dataset.org$apache$spark$sql$Dataset$$withPlan(Dataset.scala:3406)
  at org.apache.spark.sql.Dataset.select(Dataset.scala:1334)
  ... 49 elided

scala> df.select(max("Volume")).show()
+-----------+
|max(Volume)|
+-----------+
|  315541800|
+-----------+


scala> df.select(min("Volume")).show()
+-----------+
|min(Volume)|
+-----------+
|    3531300|
+-----------+


scala>

scala> val dfvol = maxV+minV
<console>:25: error: not found: value maxV
       val dfvol = maxV+minV
                   ^
<console>:25: error: not found: value minV
```

### 11-.With Scala/Spark $ resolve the next

### A-.With Scala/Spark $ resolve the next
Need to know the data less than numbesr 600 and cout that

```
df.filter($"Close"<600).count()
```

```
+----------+
|   3531300|
+----------+


scala>

scala> val dfvol = maxV+minV
<console>:25: error: not found: value maxV
       val dfvol = maxV+minV
                   ^
<console>:25: error: not found: value minV
       val dfvol = maxV+minV
                        ^

scala> val maxV =df.select(max("Volume"))
maxV: org.apache.spark.sql.DataFrame = [max(Volume): int]

scala> val minV =df.select(min("Volume"))
minV: org.apache.spark.sql.DataFrame = [min(Volume): int]

scala> val dfvol = maxV+minV
<console>:29: error: type mismatch;
 found   : org.apache.spark.sql.DataFrame
    (which expands to)  org.apache.spark.sql.Dataset[org.apache.spark.sql.Row]
 required: String
       val dfvol = maxV+minV
                       ^

scala>

scala> df.filter($"Close"<600).count()
res31: Long = 1218

scala>
```