# Experiment 7: Spark User Defined Function
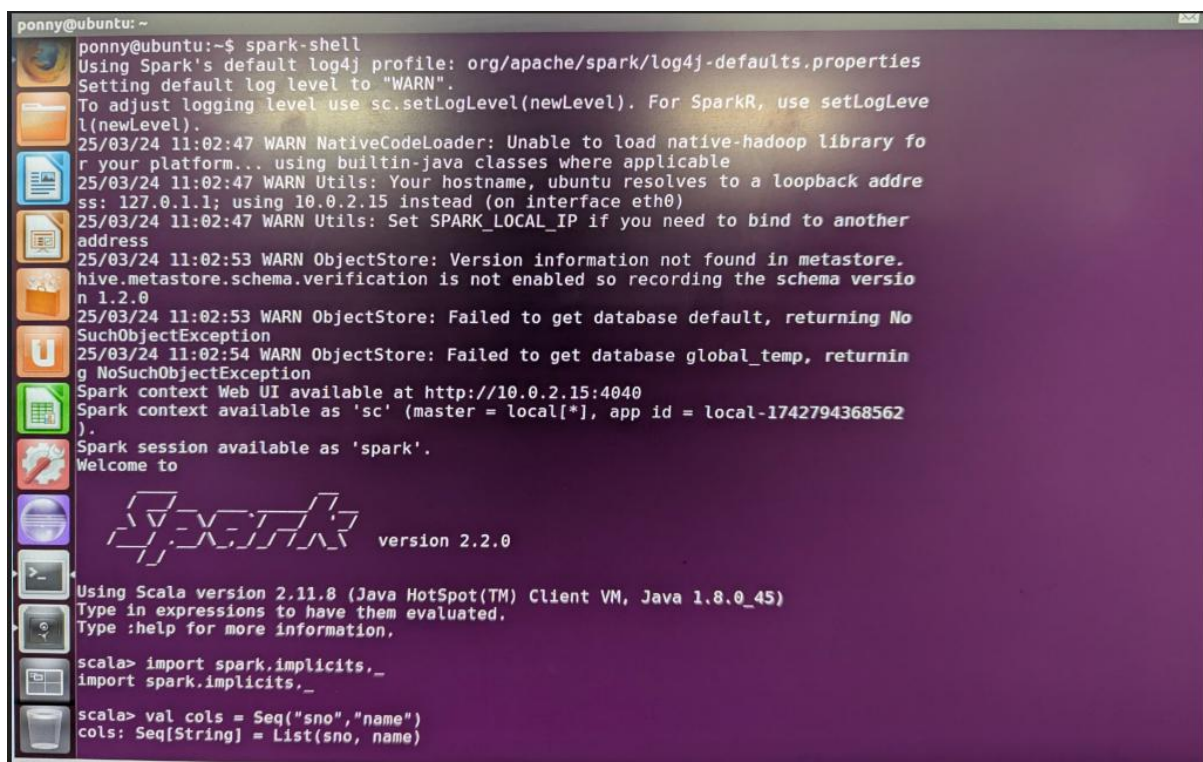
**Namansh Singh Maurya**
**22MIA1034**

## Aim:

To use User Defined Functions in Spark.

## Algorithm/Procedure:

1. Checking if Spark is present in the machine or not by using **spark-shell** command if it is present, we will spark version displayed on the screen, otherwise error.
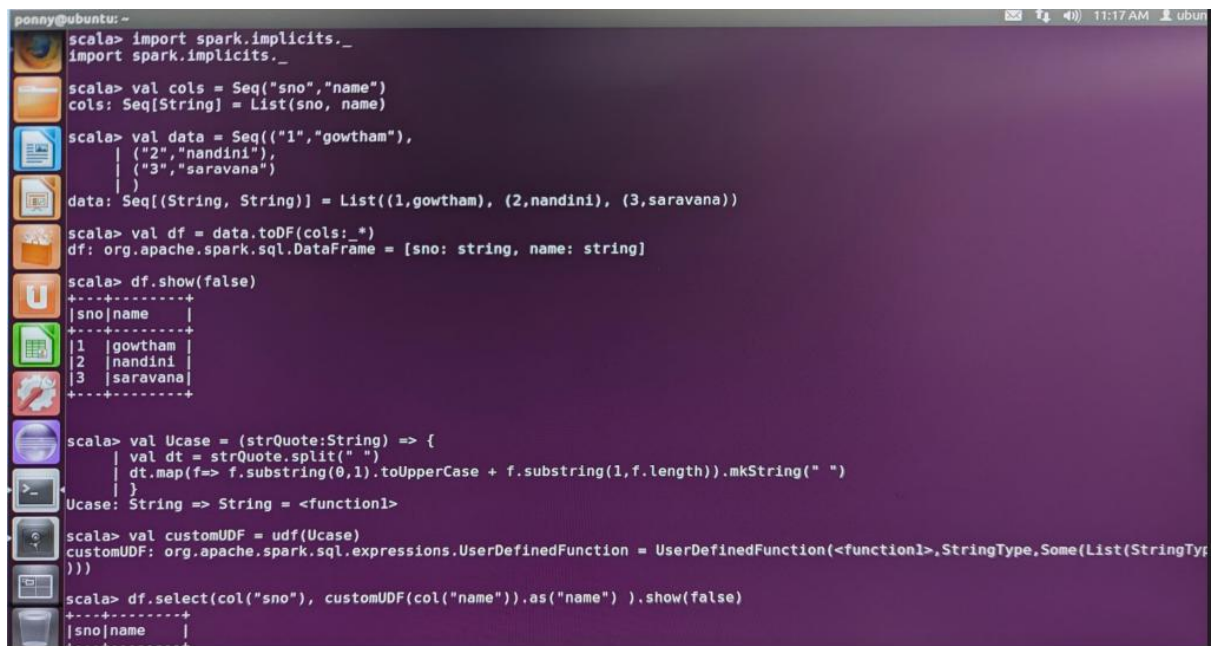


2. After Reviewing the version of Spark, we can move on to the predefined functions.

3. Below are all the codes that were used in doing this experiment.

## Program:

```
import spark.implicits._
val cols = Seq("sno","name")
val data = Seq(("1",gowtham),("2","nandini"),("3","saravana"))
val df = data.toDF(cols:_*)
df.show(false)
val Ucase  =  (strQuote: String) => { val dt = strQuote. split (" ")
dt.map(f=> f. substring (0, 1) . toUpperCase + f.substring (1, f.length) )
.mkString (" ") }
val customUDF = udf(Ucase)
df.select(col("sno"),customUDF(col("name")).as("name")).show(false)
```

## Output:

```
                | ("2","nandini"),
                | ("3","saravana")
                | )
data: Seq[(String, String)] = List((1,gowtham), (2,nandini), (3,saravana))

scala> val df = data.toDF(cols:_*)
df: org.apache.spark.sql.DataFrame = [sno: string, name: string]

scala> df.show(false)
+---+--------+
|sno|name    |
+---+--------+
|1  |gowtham |
|2  |nandini |
|3  |saravana|
+---+--------+

scala> val Ucase = (strQuote:String) => {
                | val dt = strQuote.split(" ")
                | dt.map(f=> f.substring(0,1).toUpperCase + f.substring(1,f.length)).mkString(" ")
                | }
Ucase: String => String = <function1>

scala> val customUDF = udf(Ucase)
customUDF: org.apache.spark.sql.expressions.UserDefinedFunction = UserDefinedFunction(<function1>,StringType,Some(List(StringType
)))

scala> df.select(col("sno"), customUDF(col("name")).as("name") ).show(false)
+---+--------+
|sno|name    |
+---+--------+
|1  |Gowtham |
|2  |Nandini |
|3  |Saravana|
+---+--------+

scala>
```

## Result:

Hence we used spark to use the user defined functions and get the correct output as intended.
Spark's functions were successfully understood.