# TASK1
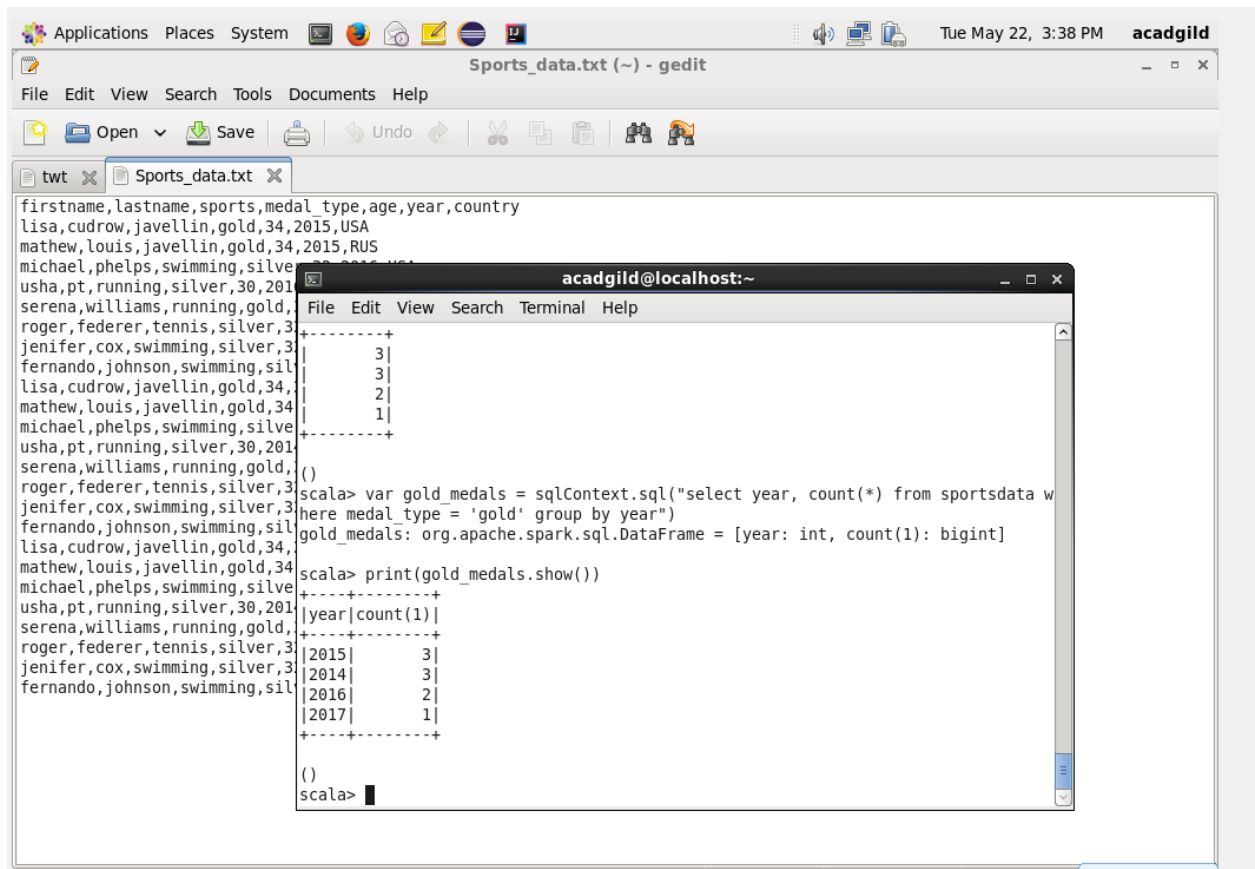
1) Spark-shell –packages com.databricks:spark-csv_2.11:1.3.0
2) Val sqlContext = new org.apache.spark.sql.SQLContext(sc)
3) Val sports_data = sqlContext.read.format("com.databricks.spark.csv").option("header", "true").option("inferSchema", "true").load("sports_data.txt")
4) Sports_data.registerTempTable("sportsdata")
5) Sports_data.printSchema()
6) Var gold_medals = sqlContext.sql("select year, count(*) from sportsdata where medal_type = 'gold' group by year")



How many silver medals have been won by the USA in each sport?
Var silver_usa = sqlContext.sql("select sports, count(*) from sportsdata where medal_type = 'silver' and country = 'USA' group by sports)
Print(silver_usa.show())

## TASK 2

1) sqlContext.udf.register("createnewname", (colA: String, colB: String) { "Mr." + colA.substring(0,2) + colB })

2) sqlContext.sql("select firstname, lastname, createnewname(firstname, lastname) from sportsdata").show()

Add new column with category names for athletes based on medals earned and age.

1) Val sqlContext = new org.apache.spark.sql.SQLContext(sc)
2) Val data = sc.textFile("/home/acadgild/Sports_data.txt")
3) Val header = data.first()
4) Data1 = data.filter(row=>row != header)
5) Case class sports_class(firstname:String, sports:String, medal_type: String, age: Int, year: Int, country: String)
6) Val sports = data1.map(x=>x.split(",")).map(x=>sports_class(x(0), x(1), x(2), x(3), x(4).toInt, x(5).toInt, x(6))).toDF
7) Sports.registerTempTable("Sports")
8) Val sports1 = sqlContext.sql("select *, IF((medal_type = 'gold' and age >= 32), 'pro', IF((medal_type = 'gold' and age <= 31), 'amateur', IF((medal_type = 'silver' and age >= 32), 'expert', IF(medal_type = 'silver' and age <= 31), 'rookie', 'none')))) as category from sports")
9) Sports1.registerTempTable("sports1")

10) sqlContext.sql("select firstname, category from sports1").show()

```
... 49 elided

scala> val sports1 = sqlContext.sql("select *, (IF((medal_type = 'gold' and age >= 32),'pro',IF((medal_type = 'gold' and age
<= 31), 'amateur', IF((medal_type = 'silver' and age >= 32), 'expert', IF((medal_type = 'silver' and age <= 31), 'rookie','n
ne')))))) as category from sports")
sports1: org.apache.spark.sql.DataFrame = [firstname: string, lastname: string ... 6 more fields]

scala> sports1.registerTempTable("sports1")
warning: there was one deprecation warning; re-run with -deprecation for details

scala> sqlContext.sql("select firstname, category from sports1").show()
+---------+--------+
|firstname|category|
+---------+--------+
|     lisa|     pro|
|   mathew|     pro|
|  michael|  expert|
|     usha|  rookie|
|   serena| amateur|
|    roger|  expert|
|  jenifer|  expert|
| fernando|  expert|
|     lisa|     pro|
|   mathew|     pro|
|  michael|  expert|
|     usha|  rookie|
|   serena| amateur|
|    roger|  expert|
|  jenifer|  expert|
| fernando|  expert|
|     lisa|     pro|
|   mathew|     pro|
|  michael|  expert|
|     usha|  rookie|
+---------+--------+
only showing top 20 rows


scala>
```