# Assignment 21

```scala
import org.apache.spark.sql.SparkSession
import org.apache.spark.sql.functions._
object Assignment_21 {

//declaring case class Sports
  case class Sports
(f_name:String,l_name:String,sports:String,medal_type:String,age:Int,year:Int,country:String
)

//creating spark session object
  def main(args: Array[String]): Unit = {
   println("hey scala")
   val spark = SparkSession
    .builder()
    .master("local")
    .appName("Spark SQL Assignment 21")
    .config("spark.some.config.option", "some-value")
    .getOrCreate()

   println("Spark Session Object created")


//loading of data from the text file
   val data =
spark.sparkContext.textFile("E:\\Avani\\Acadgild\\Assignment_21\\Sports_data.txt");
   val header=data.first()
   val data1 = data.filter(x => x != header)
    val num=println("Sports_Data->>"+data1.count())
   println("removed header")


//convert RDD into DataFrame, 's'
   import spark.implicits._
   val s =data1.map(x=>x.split(",")).map(x =>
Sports(x(0),x(1),x(2),x(3),x(4).trim.toInt,x(5).trim.toInt,x(6)))
     .toDF()
   s.show(25)
   println("Sports data")


//creating a temporary table
   s.registerTempTable("sport")
   println("temp table created")
```

## Task 1

## Using spark-sql, Find:

    a. What are the total number of gold medal winners every year?

```
val gold=spark.sql("SELECT year,count(*) from sport where medal_type='gold' group by year")
gold.show
println("Task1.1 output")
```

    b. How many silver medals have been won by USA in each sport?

```
val silver=spark.sql("SELECT sports,count(*) from sport where medal_type='silver' and country='USA' group by sports")

 silver.show

println("Task1.2 output")
```

## Task 2

## Using udfs on dataframe

    a. Change firstname, lastname columns into

Mr.first_two_letters_of_firstname<space>lastname

for example - michael, phelps becomes Mr.mi phelps

```
//creating a nameudf method

def nameudf = ((f_name:String,
l_name:String)=>"Mr.".concat(f_name.substring(0,2)).concat("").concat(l_name:String))

//registering the method as a udf and storing it into variable, fullname
val fullname= spark.sqlContext.udf.register("fullname",nameudf)

//sql query to fetch the desired output
val name = spark.sql("SELECT fullname(f_name, l_name) FROM sport")
name.show(25)
println("Task 2.1 output")
```

b. <u>Add a new column called ranking using udfs on dataframe, where :</u>

<u>gold medalist, with age >= 32 are ranked as pro</u>

<u>gold medalists, with age <= 31 are ranked amateur</u>

<u>silver medalist, with age >= 32 are ranked as expert</u>

<u>silver medalists, with age <= 31 are ranked rookie</u>

```scala
//declaring a Ranking udf along with the mentioned conditions
def Ranking = udf((medal_type: String, age: Int) => (medal_type,age) match
  {
    case (medal_type,age) if medal_type == "gold" && age >= 32 => "Pro"
    case (medal_type,age) if medal_type == "gold" && age <= 32 => "amateur"
    case (medal_type,age) if medal_type == "silver" && age >= 32 => "expert"
    case (medal_type,age) if medal_type == "silver" && age <= 32 => "rookie"
  })
//adding it (along with the parameters) as a column to the 's' dataframe, declared above
  val r = s.withColumn("Rank", Ranking(s("medal_type"),s("age")))
r.show(25)
  println("Task 2.2 output")
```
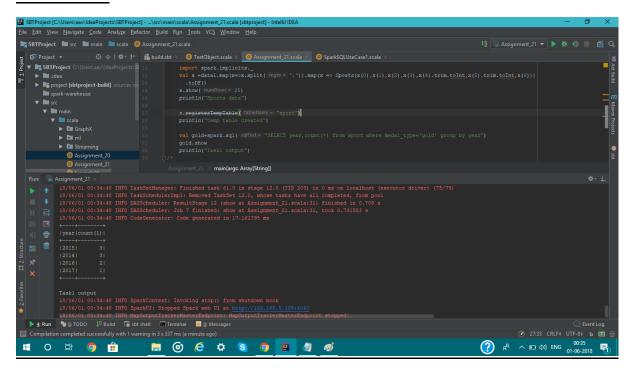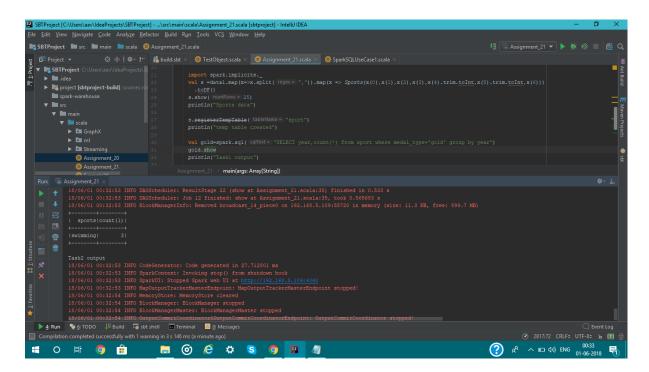
## Complete code

```
import org.apache.spark.sql.SparkSession
import org.apache.spark.sql.functions._
object Assignment_21 {
  case class Sports
(f_name:String,l_name:String,sports:String,medal_type:String,age:Int,year:Int,country:String
)
  def main(args: Array[String]): Unit = {
    println("hey scala")
    val spark = SparkSession
      .builder()
      .master("local")
      .appName("Spark SQL Assignment 21")
      .config("spark.some.config.option", "some-value")
      .getOrCreate()

    println("Spark Session Object created")

    val data =
spark.sparkContext.textFile("E:\\Avani\\Acadgild\\Assignment_21\\Sports_data.txt");
    val header=data.first()
    val data1 = data.filter(x => x != header)
    val num=println("Sports_Data->>"+data1.count())
    println("removed header")

    import spark.implicits._
    val s =data1.map(x=>x.split(",")).map(x =>
Sports(x(0),x(1),x(2),x(3),x(4).trim.toInt,x(5).trim.toInt,x(6)))
      .toDF()
    s.show(25)
    println("Sports data")

    s.registerTempTable("sport")
    println("temp table created")

    val gold=spark.sql("SELECT year,count(*) from sport where medal_type='gold' group by
year")
    gold.show
    println("Task1.1 output")


    val silver=spark.sql("SELECT sports,count(*) from sport where medal_type='silver' and
country='USA' group by sports")
```
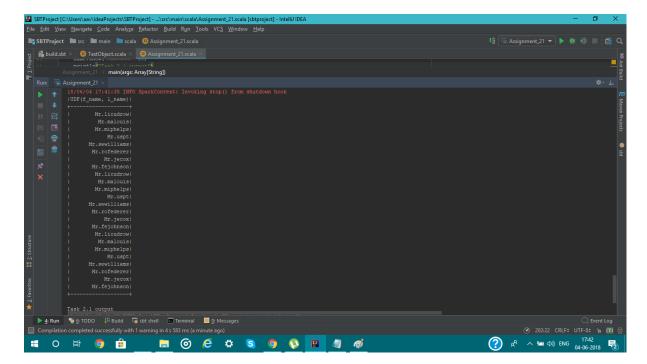
```scala
    silver.show
    println("Task1.2 output")

  def nameudf = ((f_name:String,
l_name:String)=>"Mr.".concat(f_name.substring(0,2)).concat("").concat(l_name:String))
    val fullname= spark.sqlContext.udf.register("fullname",nameudf)
    val name = spark.sql("SELECT fullname(f_name, l_name) FROM sport")
    name.show(25)
    println("Task 2.1 output")


    def Ranking = udf((medal_type: String, age: Int) => (medal_type,age) match
    {
      case (medal_type,age) if medal_type == "gold" && age >= 32 => "Pro"
      case (medal_type,age) if medal_type == "gold" && age <= 32 => "amateur"
      case (medal_type,age) if medal_type == "silver" && age >= 32 => "expert"
      case (medal_type,age) if medal_type == "silver" && age <= 32 => "rookie"
    })
 // val RANK= spark.sqlContext.udf.register("RANK",Ranking)
    val r = s.withColumn("Rank", Ranking(s("medal_type"),s("age")))
r.show(25)
    println("Task 2.2 output")

 }

}
```
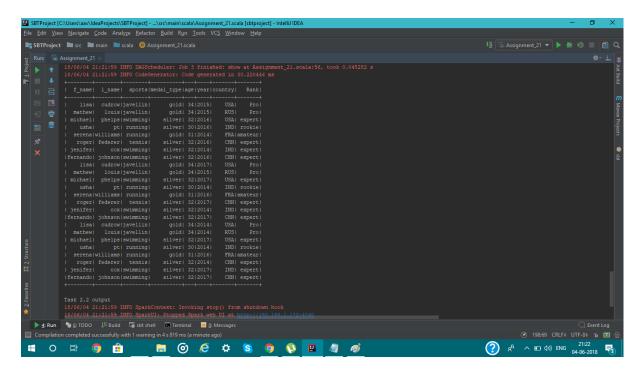
## Screenshots



**1.a.**



**1.b.**

**2.a.**



**2.b.**