

Assignment 20

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

object Assignment_20 {

  //Create a case class globally to be used inside the main method
  case class Transport(mode: String, amount: Int)
  case class User(id:Int, name:String,age :Int)
  case class Holiday (t_id:Int,Source:String,Destination: String,mode: String, distance:Int,
year:Int)

  def main(args: Array[String]): Unit = {

    //Create a spark session object
    val spark = SparkSession
      .builder()
      .master("local")
      .appName("Spark SQL basic example")
      .config("spark.some.config.option", "some-value")
      .getOrCreate()

    println("Spark Session Object created")

    import spark.implicits._          //to convert RDD into DataFrame
    val hl = spark.sparkContext
      .textFile("E:\\Avani\\Acadgild\\Assignment_20\\Data\\Holidays.txt")
      .map(_._split(","))
      .map(attributes => Holiday(attributes(0).trim.toInt,
attributes(1),attributes(2),attributes(3),attributes(4).trim.toInt,attributes(5).trim.toInt))
      .toDF()

    hl.show()

    println("Holiday data")
  }
}
```

```

val trans = spark.sparkContext
    .textFile("E:\\Avani\\Acadgild\\Assignment_20\\Data\\Transport.txt")
    .map(_ .split(", "))
    .map(attributes => Transport(attributes(0), attributes(1).trim.toInt))
    .toDF()
trans.show()

println("Transport data")

```

```

val user = spark.sparkContext
    .textFile("E:\\Avani\\Acadgild\\Assignment_20\\Data\\User.txt")
    .map(_ .split(", "))
    .map(attributes => User(attributes(0).trim.toInt, attributes(1), attributes(2).trim.toInt))
    .toDF()
user.show()

println("User data")

//Creating temporary tables
hl.registerTempTable("holiday")
user.registerTempTable("people")
trans.registerTempTable("transport")

```

a. What is the distribution of the total number of air-travellers per year?

```

val total = spark.sql("SELECT year,COUNT(*) from holiday group by year")
total.show()
println("Task1 output")

```

b. What is the total air distance covered by each user per year?

```

val dist = spark.sql("SELECT t_id,year,sum(distance) from holiday group by t_id,year
order by t_id,year")

```

```
dist.show()
println("Task2 output")
```

c. Which user has travelled the largest distance till date?

```
val max_d = spark.sql("SELECT t_id,max(d) from (SELECT t_id, sum(distance) d from
holiday group by t_id) group by t_id")
max_d.show(1)
println("Task3 output")
```

d. What is the most preferred destination for all users?

```
val ct=spark.sql("SELECT Destination from holiday group by 1 having count
(Destination)= (select max(c) from (SELECT Destination,count(Destination) c from
holiday group by Destination))")

ct.show()

println("Task4 output")
```

e. Which route is generating the most revenue per year?

```
val rev = spark.sql("SELECT h.source, h.destination, sum(h.distance*t.amount) revenue
from holiday h join transport t ON h.mode=t.mode group by h.source,h.destination")

val x= rev.sort($"revenue".desc).first

println("most revenue generated route-->" + x)
```

f. What is the total amount spent by every user on air-travel per year?

```
val amt = spark.sql("SELECT h.t_id, h.year,sum(t.amount) from holiday h join transport t
ON h.mode=t.mode group by h.t_id,h.year")

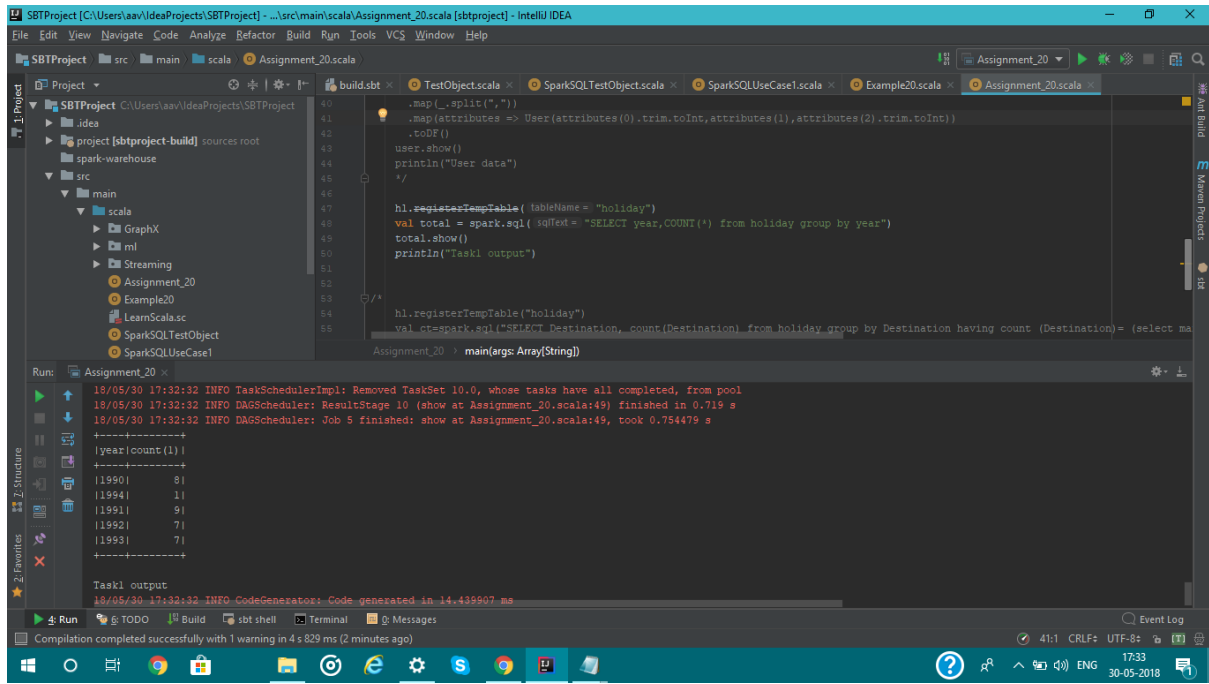
amt.show(32)

println("Task6")
```

- g. Considering age groups of < 20 , 20-35, 35 > ,Which age group is travelling the most every year

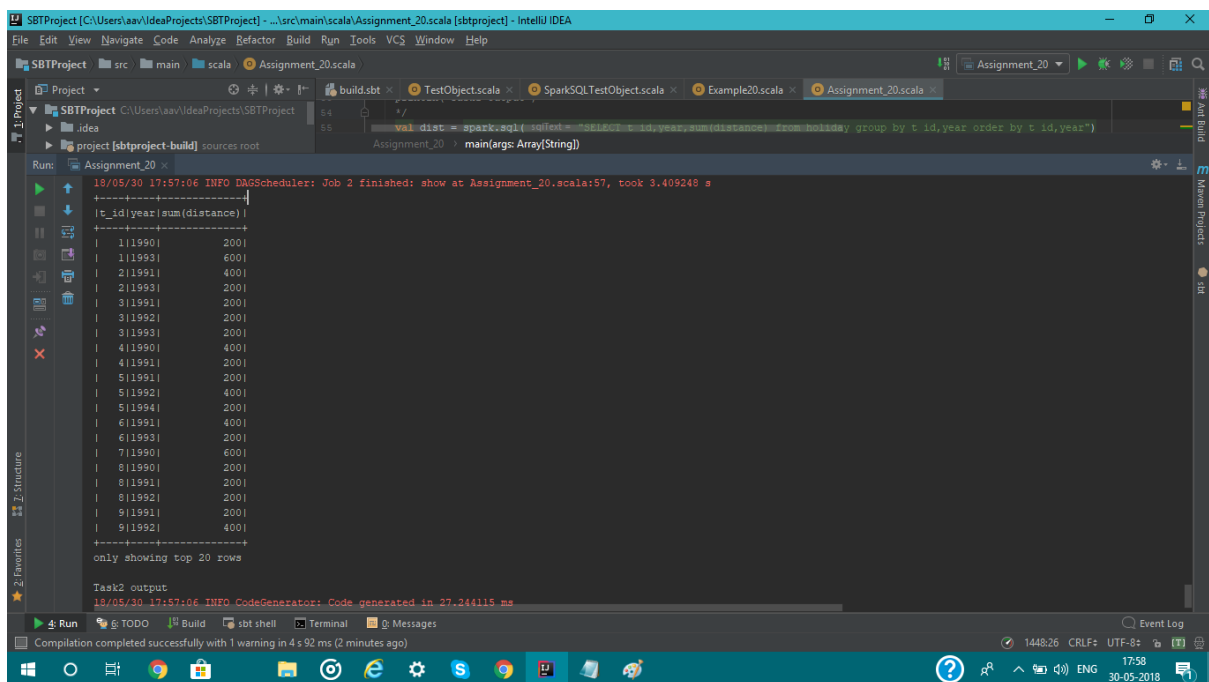
```
val new_user = spark.sql("SELECT * from holiday h join people p ON h.t_id = p.id")
new_user.show()
new_user.registerTempTable("NewDetails")
val age = spark.sql("select year, age ,count(age) age_cnt from NewDetails group by
year,age order by age_cnt desc").take(1)
println("Task7 output--->" + age.foreach(println))
```

Screenshot



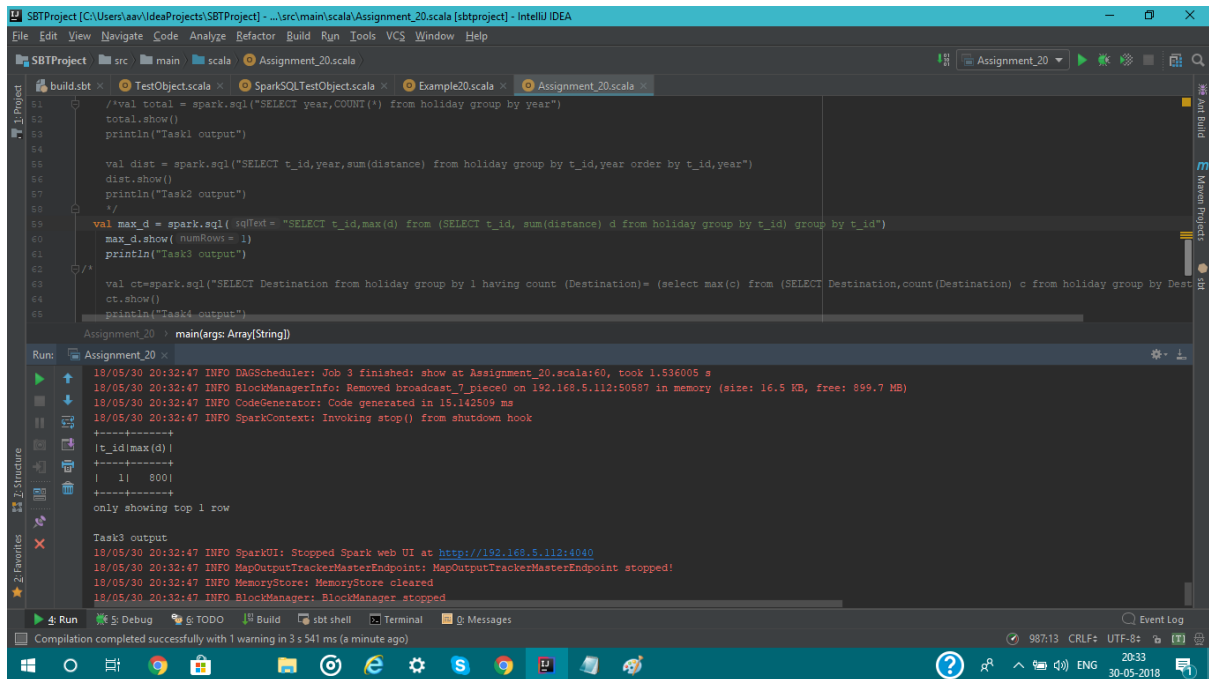
```
File Edit View Navigate Code Analyze Refactor Build Run Tools VCS Window Help
SBTProject [C:\Users\aa\IdeaProjects\SBTProject] - \src\main\scala\Assignment_20.scala [sbtproject] - IntelliJ IDEA
Project: SBTProject
  - src
    - main
      - scala
        - Assignment_20.scala
Run: Assignment_20
18/05/30 17:32:32 INFO TaskSchedulerImpl: Removed TaskSet 10.0, whose tasks have all completed, from pool
18/05/30 17:32:32 INFO DAGScheduler: ResultStage 10 (show at Assignment_20.scala:49) finished in 0.719 s
18/05/30 17:32:32 INFO DAGScheduler: Job 5 finished: show at Assignment_20.scala:49, took 0.754479 s
+-----+
|year|count|
+-----+
|1990| 8|
|1991| 1|
|1992| 9|
|1993| 7|
|1994| 7|
+-----+
Task1 output
18/05/30 17:32:32 INFO CodeGenerator: Code generated in 14.439507 ms
Compilation completed successfully with 1 warning in 4 s 829 ms (2 minutes ago)
```

a.



```
File Edit View Navigate Code Analyze Refactor Build Run Tools VCS Window Help
SBTProject [C:\Users\aa\IdeaProjects\SBTProject] - \src\main\scala\Assignment_20.scala [sbtproject] - IntelliJ IDEA
Project: SBTProject
  - src
    - main
      - scala
        - Assignment_20.scala
Run: Assignment_20
18/05/30 17:57:06 INFO DAGScheduler: Job 2 finished: show at Assignment_20.scala:57, took 3.409248 s
+-----+
|t_id|year|sum(distance)|
+-----+
|1|1990| 200|
|1|1993| 600|
|2|1991| 400|
|2|1993| 200|
|3|1991| 200|
|3|1992| 200|
|3|1993| 200|
|4|1990| 400|
|4|1991| 200|
|5|1991| 200|
|5|1992| 400|
|5|1994| 200|
|6|1991| 400|
|6|1993| 200|
|7|1990| 600|
|8|1990| 200|
|8|1991| 200|
|8|1992| 200|
|9|1991| 200|
|9|1992| 400|
+-----+
only showing top 20 rows
Task2 output
18/05/30 17:57:06 INFO CodeGenerator: Code generated in 27.244115 ms
Compilation completed successfully with 1 warning in 4 s 92 ms (2 minutes ago)
```

b.



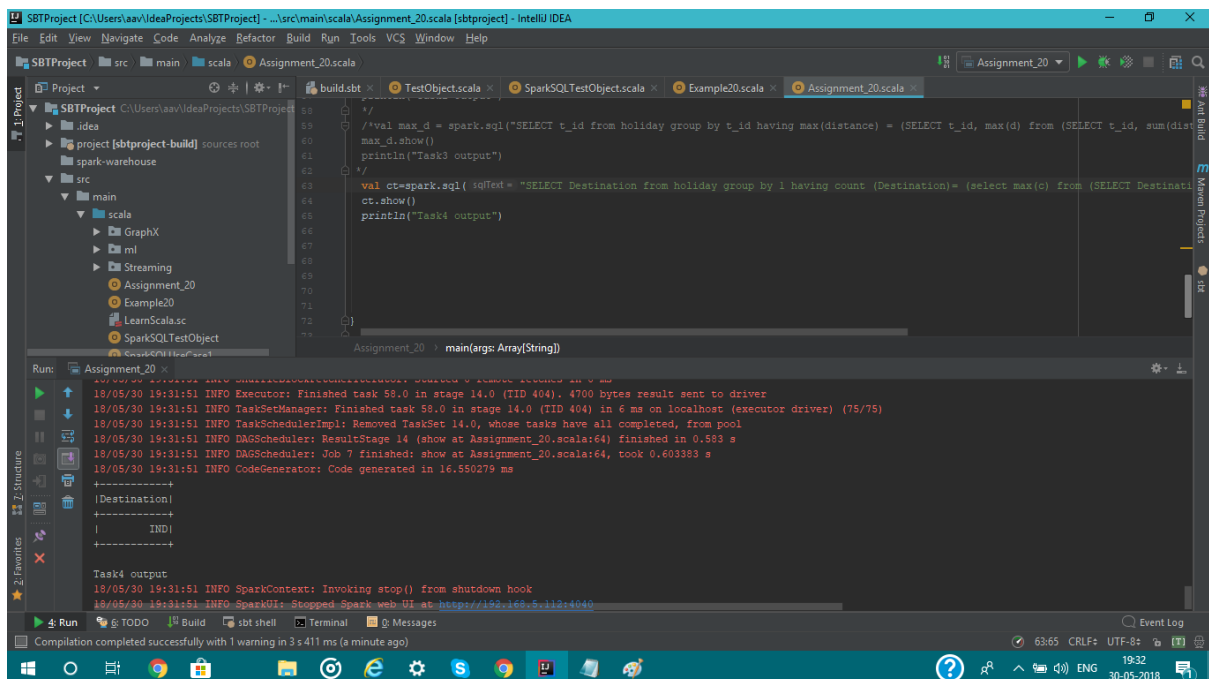
```
117 //val total = spark.sql("SELECT year,COUNT(*) from holiday group by year")
118 total.show()
119 println("Task1 output")
120
121 val dist = spark.sql("SELECT t_id,year,sum(distance) from holiday group by t_id,year order by t_id,year")
122 dist.show()
123 println("Task2 output")
124
125 val max_d = spark.sql( sqlText = "SELECT t_id,max(d) from (SELECT t_id, sum(distance) d from holiday group by t_id) group by t_id")
126 max_d.show( numRows = 1)
127 println("Task3 output")
128
129 val ct=spark.sql("SELECT Destination from holiday group by l having count (Destination)= (select max(c) from (SELECT Destination,count(Destination) c from holiday group by Destination) group by Destination)")
130 ct.show()
131 println("Task4 output")
132
133 Assignment_20 main(args: Array[String])
```

Run Assignment_20

```
18/05/30 20:32:47 INFO DAGScheduler: Job 3 finished: show at Assignment_20.scala:60, took 1.536005 s
18/05/30 20:32:47 INFO BlockManagerInfo: Removed broadcast_7_piece0 on 192.168.5.112:50507 in memory (size: 16.5 KB, free: 899.7 MB)
18/05/30 20:32:47 INFO CodeGenerator: Code generated in 15.142509 ms
18/05/30 20:32:47 INFO SparkContext: Invoking stop() from shutdown hook
+-----+
|t_id|max(d)|
+-----+
| 1 | 800 |
+-----+
only showing top 1 row

Task3 output
18/05/30 20:32:47 INFO SparkUI: Stopped Spark web UI at http://192.168.5.112:4040
18/05/30 20:32:47 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
18/05/30 20:32:47 INFO MemoryStore: MemoryStore cleared
18/05/30 20:32:47 INFO BlockManager: BlockManager stopped
```

c.



```
134 //val max_d = spark.sql("SELECT t_id,max(d) from (SELECT t_id, sum(distance) d from holiday group by t_id) group by t_id")
135 max_d.show()
136 println("Task3 output")
137
138 val ct=spark.sql( sqlText = "SELECT Destination from holiday group by l having count (Destination)= (select max(c) from (SELECT Destination,count(Destination) c from holiday group by Destination) group by Destination)")
139 ct.show()
140 println("Task4 output")
141
142 Assignment_20 main(args: Array[String])
```

Run Assignment_20

```
18/05/30 19:31:51 INFO DAGScheduler: Job 7 finished: show at Assignment_20.scala:64, took 0.603383 s
18/05/30 19:31:51 INFO CodeGenerator: Code generated in 16.550279 ms
+-----+
|Destination|
+-----+
| IND |
+-----+

Task4 output
18/05/30 19:31:51 INFO SparkContext: Invoking stop() from shutdown hook
18/05/30 19:31:51 INFO SparkUI: Stopped Spark web UI at http://192.168.5.112:4040
```

d.

```
164 ct.show()
165 println("Task4 output")
166
167 val rev = spark.sql("SELECT h.source, h.destination, sum(h.distance*t.amount) revenue from holiday h join transport t ON h.mode=t.mode group by h.source,h.destination")
168 val x= rev.sort($"revenue".desc).first
169 println("most revenue generated route-->" + x)
170
171
172
```

Run: Assignment_20 x

```
18/05/30 22:36:18 INFO ShuffleBlockFetcherIterator: Started 0 remote fetches in 0 ms
18/05/30 22:36:18 INFO Executor: Finished task 195.0 in stage 3.0 (TID 401). 5333 bytes result sent to driver
18/05/30 22:36:18 INFO TaskSetManager: Finished task 195.0 in stage 3.0 (TID 401) in 15 ms on localhost (executor driver) (200/200)
18/05/30 22:36:18 INFO TaskSchedulerImpl: Removed TaskSet 3.0, whose tasks have all completed, from pool
18/05/30 22:36:18 INFO DAGScheduler: ResultStage 3 (first at Assignment_20.scala:69) finished in 1.262 s
18/05/30 22:36:18 INFO DAGScheduler: Job 0 finished: first at Assignment_20.scala:69, took 6.926237 s
18/05/30 22:36:18 INFO CodeGenerator: Code generated in 16.330502 ms
most revenue generated route-->[CHN,IND,136000]
18/05/30 22:36:18 INFO SparkContext: Invoking stop() from shutdown hook
18/05/30 22:36:18 INFO SparkUI: Stopped Spark web UI at http://159.254.29.30:4040
18/05/30 22:36:18 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
18/05/30 22:36:18 INFO MemoryStore: MemoryStore cleared
18/05/30 22:36:18 INFO BlockManager: BlockManager stopped
18/05/30 22:36:18 INFO BlockManagerMaster: BlockManagerMaster stopped
18/05/30 22:36:18 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
18/05/30 22:36:18 INFO SparkContext: Successfully stopped SparkContext
18/05/30 22:36:18 INFO ShutdownHookManager: Shutdown hook called
18/05/30 22:36:18 INFO ShutdownHookManager: Deleting directory C:\Users\aa\AppData\Local\Temp\spark-0568e07d-c4e4-47a2-b152-122e38c999d6

Process finished with exit code 0
```

Compilation completed successfully with 1 warning in 3 s 579 ms (4 minutes ago)

e.

```
76 rev_user.show()
77
```

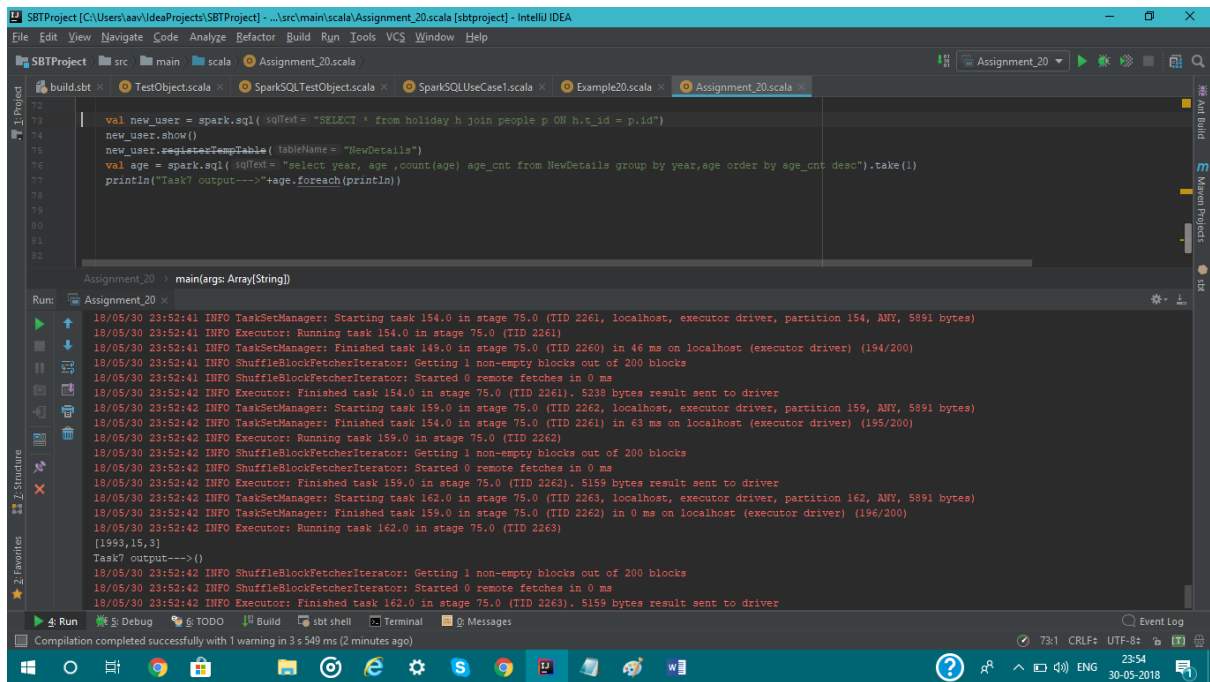
Run: Assignment_20 x

```
-----+-----+
|t_id|year|sum(amount)|
-----+-----+
| 3|1991|      170|
| 6|1993|      170|
| 3|1992|      170|
| 7|1990|      510|
|10|1993|      170|
| 6|1991|      340|
| 2|1991|      340|
| 4|1991|      170|
| 5|1991|      170|
| 5|1994|      170|
| 8|1991|      170|
|11|1990|      170|
| 5|1992|      340|
| 4|1990|      340|
| 3|1993|      170|
|10|1990|      170|
| 2|1993|      170|
|11|1993|      510|
| 9|1991|      170|
| 9|1992|      340|
| 8|1990|      170|
|10|1992|      170|
| 8|1992|      170|
-----+-----+

Task6 output
18/05/30 23:57:29 INFO SparkUI: Stopped Spark web UI at http://192.168.5.112:4040
```

Compilation completed successfully with 1 warning in 3 s 901 ms (a minute ago)

f.



g.