# ASSIGNMENT 19.1

## Input Dataset:

We have sports related data with us:

**Sports_data**

firstname, lastname, sports, medal_type, age, year, country

```
  Sports_data.txt ☒
   1   firstname,lastname,sports,medal_type,age,year,country
   2   lisa,cudrow,javellin,gold,34,2015,USA
   3   mathew,louis,javellin,gold,34,2015,RUS
   4   michael,phelps,swimming,silver,32,2016,USA
   5   usha,pt,running,silver,30,2016,IND
   6   serena,williams,running,gold,31,2014,FRA
   7   roger,federer,tennis,silver,32,2016,CHN
   8   jenifer,cox,swimming,silver,32,2014,IND
   9   fernando,johnson,swimming,silver,32,2016,CHN
  10   lisa,cudrow,javellin,gold,34,2017,USA
  11   mathew,louis,javellin,gold,34,2015,RUS
  12   michael,phelps,swimming,silver,32,2017,USA
  13   usha,pt,running,silver,30,2014,IND
  14   serena,williams,running,gold,31,2016,FRA
  15   roger,federer,tennis,silver,32,2017,CHN
  16   jenifer,cox,swimming,silver,32,2014,IND
  17   fernando,johnson,swimming,silver,32,2017,CHN
  18   lisa,cudrow,javellin,gold,34,2014,USA
  19   mathew,louis,javellin,gold,34,2014,RUS
  20   michael,phelps,swimming,silver,32,2017,USA
  21   usha,pt,running,silver,30,2014,IND
  22   serena,williams,running,gold,31,2016,FRA
  23   roger,federer,tennis,silver,32,2014,CHN
  24   jenifer,cox,swimming,silver,32,2017,IND
  25   fernando,johnson,swimming,silver,32,2017,CHN
```

## Problem Statement:

Using spark-SQL, find:

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

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

## Solution:

**1.** Here is the Spark code snippet to find the total number of gold medal winners every year:

```
// import required Spark packages

import org.apache.spark.sql.SparkSession

object Assignment18_1 {
  def main(args: Array[String]): Unit = {
// create a SparkSession object that can be used to create various contexts of Spark such as sqlContext
    val spark = SparkSession
      .builder()
      .config("spark.sql.warehouse.dir", "file:///c:/tmp/spark-warehouse")
      .master("local[*]")
      .getOrCreate()

// initialize sqlContext

    val sqlContext = spark.sqlContext

// load input data file – Sports_data.txt

    val sports_DF = sqlContext.read.option("header", "true").csv("E:\\Acadgild\\Session 19\\Sports_data.txt")

//create a temporary view – Sports_data

    sports_DF.createOrReplaceTempView("Sports_data")

// SQL query to find the total number of gold medal winners every year
    sqlContext.sql("SELECT year AS Year, COUNT(*) AS Gold_Medals_Won "+
      "FROM Sports_data " +
      "WHERE medal_type = 'gold' " +
      "GROUP BY year ").show()
  }
}
```

```scala
5  ▶   object Assignment19_1 {
6
7  ▶     def main(args: Array[String]): Unit = {
8
9          val spark = SparkSession
10           .builder()
11           .config("spark.sql.warehouse.dir", "file:///c:/tmp/spark-warehouse")
12           .master("local[*]")
13           .getOrCreate()
14         val sqlContext = spark.sqlContext
15         val sports_DF = sqlContext.read.option("header", "true").csv("E:\\Acadgild\\Session 19\\Sports_data.txt")
16         sports_DF.createOrReplaceTempView("Sports_data")
17
18         sqlContext.sql("SELECT year AS Year, COUNT(*) AS Gold_Medals_Won "+
19             "FROM Sports_data " +
20             "WHERE medal_type = 'gold'" +
21             "GROUP BY year ").show()
22       }
23   }
```

**Output:**

```
+----+---------------+
|Year|Gold_Medals_Won|
+----+---------------+
|2016|              2|
|2017|              1|
|2014|              3|
|2015|              3|
+----+---------------+
```

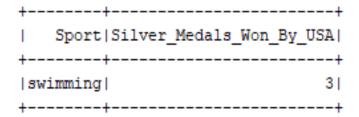2. Here is the Spark code to find the total number of silver medals won by USA in each sport:

// import required Spark packages

import org.apache.spark.sql.SparkSession

object Assignment18_1 {

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

// create a SparkSession object that can be used to create various contexts of Spark such as sqlContext

    val spark = SparkSession

```
        .builder()
        .config("spark.sql.warehouse.dir", "file:///c:/tmp/spark-warehouse")
        .master("local[*]")
        .getOrCreate()
```

// initialize sqlContext

```
    val sqlContext = spark.sqlContext
```

// load input data file – Sports_data.txt

```
    val sports_DF = sqlContext.read.option("header", "true").csv("E:\\Acadgild\\Session 19\\Sports_data.txt")
```

//create a temporary view – Sports_data

```
    sports_DF.createOrReplaceTempView("Sports_data")
```

// SQL query to find the total number of gold medal winners every year
```
    sqlContext.sql("SELECT sports AS Sport, COUNT(*) AS Silver_Medals_Won_By_USA "+
        "FROM Sports_data " +
        "WHERE medal_type = 'silver' AND country = 'USA' " +
        "GROUP BY sports ").show()
  }
}
```



```scala
 4
 5  ▶   ⊟object Assignment19_1 {
 6
 7  ▶   ⊟   def main(args: Array[String]): Unit = {
 8
 9          val spark = SparkSession
10            .builder()
11            .config("spark.sql.warehouse.dir", "file:///c:/tmp/spark-warehouse")
12            .master("local[*]")
13            .getOrCreate()
14          val sqlContext = spark.sqlContext
15          val sports_DF = sqlContext.read.option("header", "true").csv("E:\\Acadgild\\Session 19\\Sports_data.txt")
16          sports_DF.createOrReplaceTempView("Sports_data")
17
18          sqlContext.sql("SELECT sports AS Sport, COUNT(*) AS Silver_Medals_Won_By_USA "+
19              "FROM Sports_data " +
20              "WHERE medal_type = 'silver' AND country = 'USA' " +
21              "GROUP BY sports ").show()
22        }
23  ⊟}
```

**Output:**

```
+--------+----------------------+
|   Sport|Silver_Medals_Won_By_USA|
+--------+----------------------+
|swimming|                     3|
+--------+----------------------+
```

To summarize, USA has won **3** silver medals in swimming category.