# BIG DATA HADOOP AND SPARK DEVLOPMENT

# ASSIGNMENT 28

Table of Contents:

# BIG DATA HADOOPAND SPARK DEVELOPMENT

## 1. Introduction

In this assignment, the given tasks are performed and Output of the tasks are recorded in the form of Screenshots.

## 2. Objective

This Assignment consolidates the deeper understanding of the Session 28 Spark MLIB - I

## 3. Problem Statement

- Task 1
    - Find out the top 5 most visited destinations.
- Task 2
    - Which month has seen the most number of cancellations due to bad weather?
- Task 3
    - Which route (origin & destination) has seen the maximum diversion?

## 4. Expected Output

**Aviation data analysis**

You can download the datasets from the following links: Delayed_Flights.csv There are 29 columns in this dataset. Some of them have been mentioned below:

• Year: 1987 – 2008

• Month: 1 – 12

• FlightNum: Flight number

• Canceled: Was the flight canceled?

• CancelleationCode: The reason for cancellation. Now the very first thing is that we are going to implement this using Spark SQL. So as per requirement, we proceed to set up the Spark Context and load the input CSV file as shown below.

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

object DelayedFlightsAnalysis {

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

    println("hey scala")

    val spark = SparkSession
      .builder()
      .master( master = "local")
      .appName( name = "Delayed Flight Analysis")
      .config("spark.some.config.option", "some-value")
      .getOrCreate()

    println("Spark Session Object created")

    //Set the log level as warning
    spark.sparkContext.setLogLevel("WARN")
```

Now to load the file.

```scala
    val df1 = spark.sqlContext.read
      .option("header", "true")
      .option("inferSchema", "true")
      .csv( path = "C:\\Users\\Ankith M\\Desktop\\Hadoop\\Spark\\DelayedFlights.csv")

    println("Spark Delayed flight DF1 created!")

    df1.show()

    df1.printSchema()
```

Output of DataFrame created after reading the file and schema of the file.



DelayedFlightsAnalysis ×

```
root
 |-- _c0: integer (nullable = true)
 |-- Year: integer (nullable = true)
 |-- Month: integer (nullable = true)
 |-- DayofMonth: integer (nullable = true)
 |-- DayOfWeek: integer (nullable = true)
 |-- DepTime: double (nullable = true)
 |-- CRSDepTime: integer (nullable = true)
 |-- ArrTime: double (nullable = true)
 |-- CRSArrTime: integer (nullable = true)
 |-- UniqueCarrier: string (nullable = true)
 |-- FlightNum: integer (nullable = true)
 |-- TailNum: string (nullable = true)
 |-- ActualElapsedTime: double (nullable = true)
 |-- CRSElapsedTime: double (nullable = true)
 |-- AirTime: double (nullable = true)
 |-- ArrDelay: double (nullable = true)
 |-- DepDelay: double (nullable = true)
 |-- Origin: string (nullable = true)
 |-- Dest: string (nullable = true)
 |-- Distance: integer (nullable = true)
 |-- TaxiIn: double (nullable = true)
 |-- TaxiOut: double (nullable = true)
 |-- Cancelled: integer (nullable = true)
 |-- CancellationCode: string (nullable = true)
 |-- Diverted: integer (nullable = true)
 |-- CarrierDelay: double (nullable = true)
 |-- WeatherDelay: double (nullable = true)
 |-- NASDelay: double (nullable = true)
 |-- SecurityDelay: double (nullable = true)
 |-- LateAircraftDelay: double (nullable = true)
```
pilation completed successfully in 2s 564ms (9 minutes ago)

Now proceed to create a temporary view as below –

```
DelayedFlightsAnalysis.scala    build.sbt
32          dfl.createOrReplaceTempView( viewName = "delayed_flights")
33
34          println("temporary view for delayed flights created!!!")
35
```

Output –

```
DelayedFlightsAnalysis

 temporary view for delayed flights created!!!
```

Once the table is registered as view now we can proceed to use Spark SQL to meet each of the Problem Statements one by one.

- ## Task 1
    - ### **Find out the top 5 most visited destinations.**

```
DelayedFlightsAnalysis.scala    build.sbt
36
37          // Problem Statement 1 - Find out the top 5 most visited destinations.
38          println("the top 5 most visited destinations are: ")
39
40          val top5DF = spark.sql(
41            """select Dest, count(Dest) as Dest_Count
42              |from delayed_flights
43              |group by Dest
44              |order by Dest_Count desc
45              |limit 5
46            """.stripMargin).show()
47
```

Output:

```
DelayedFlightsAnalysis
  the top 5 most visited destinations are:
  +----+----------+
  |Dest|Dest_Count|
  +----+----------+
  | ORD|    108984|
  | ATL|    106898|
  | DFW|     70657|
  | DEN|     63003|
  | LAX|     59969|
  +----+----------+
```

- **Task 2:**
  - ○ **Which month has seen the most number of cancellations due to bad weather?**

```
DelayedFlightsAnalysis.scala    build.sbt

        // Problem Statement 2 - Which month has seen the most number of cancellations due to bad weather?
        println("the month has seen the most number of cancellations due to bad weather is: ")
    val cancelBadWeatherDF = spark.sql(
        """select Month, count(Cancelled) as Cancelled_Counts
          |from delayed_flights
          |where Cancelled = 1 and CancellationCode ='B'
          |group by Month
          |order by Cancelled_Counts desc
          |limit 1
        """.stripMargin)
    cancelBadWeatherDF.show()
```

Output:

```
DelayedFlightsAnalysis

the month has seen the most number of cancellations due to bad weather is:
+-----+----------------+
|Month|Cancelled_Counts|
+-----+----------------+
|   12|             250|
+-----+----------------+
```

- Task 3
  - **Which route (origin & destination) has seen the maximum diversion?**

```scala
//Problem Statement 3 - Which route (origin & destination) has seen the maximum diversion?
println("the route (origin & destination) has seen the maximum diversions are: ")
val diversionDF = spark.sql(
  """select Origin, Dest, count(Diverted) as Diversions_Count from delayed_flights
    |where Diverted = 1
    |group by Origin, Dest
    |order by Diversions_Count desc
    |limit 10
  """.stripMargin).show()
```

Output –

```
the route (origin & destination) has seen the maximum diversions are:
+------+----+----------------+
|Origin|Dest|Diversions_Count|
+------+----+----------------+
|   ORD| LGA|              39|
|   DAL| HOU|              35|
|   DFW| LGA|              33|
|   ATL| LGA|              32|
|   ORD| SNA|              31|
|   MIA| LGA|              31|
|   SLC| SUN|              31|
|   BUR| JFK|              29|
|   HRL| HOU|              28|
|   BUR| DFW|              25|
+------+----+----------------+


Process finished with exit code 0
```

Please find below, the complete code for this use case as a whole.

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

object DelayedFlightsAnalysis {

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

    println("hey scala")

    val spark = SparkSession
      .builder()
      .master("local")
      .appName("Delayed Flight Analysis")
      .config("spark.some.config.option", "some-value")
      .getOrCreate()

    println("Spark Session Object created")

    //Set the log level as warning
    spark.sparkContext.setLogLevel("WARN")

    val df1 = spark.sqlContext.read
      .option("header", "true")
      .option("inferSchema", "true")
```

```scala
    df1.createOrReplaceTempView("delayed_flights")

    println("temporary view for delayed flights created!!!")


    // Problem Statement 1 - Find out the top 5 most visited destinations.
    println("the top 5 most visited destinations are: ")

    val top5DF = spark.sql(
      """select Dest, count(Dest) as Dest_Count
        |from delayed_flights
        |group by Dest
        |order by Dest_Count desc
        |limit 5
      """.stripMargin).show()

    // Problem Statement 2 - Which month has seen the most number of cancellations due
to bad weather?
    println("the month has seen the most number of cancellations due to bad weather
is: ")
    val cancelBadWeatherDF = spark.sql(
      """select Month, count(Cancelled) as Cancelled_Counts
        |from delayed_flights
        |where Cancelled = 1 and CancellationCode ='B'
        |group by Month
        |order by Cancelled_Counts desc
        |limit 1
      """.stripMargin)
    cancelBadWeatherDF.show()

    //Problem Statement 3 - Which route (origin & destination) has seen the maximum
diversion?
    println("the route (origin & destination) has seen the maximum diversions are: ")
    val diversionDF = spark.sql(
      """select Origin, Dest, count(Diverted) as Diversions_Count from delayed_flights
        |where Diverted = 1
        |group by Origin, Dest
        |order by Diversions_Count desc
        |limit 10
      """.stripMargin).show()


  }
}
```