Hospital CASE STUDY

**Dataset Description**

- **DRG Definition**: The code and description identifying the MS-DRG. MS-DRGs are a classification system that groups similar clinical conditions (diagnoses) and procedures furnished by the hospital during their stay.

- **Provider Id**: The CMS Certification Number (CCN) assigned to the Medicare-certified hospital facility.

- **Provider Name**: The name of the provider.

- **Provider Street Address**: The provider's street address.

- **Provider City**: The city where the provider is located.

- **Provider State**: The state where the provider is located.

- **Provider Zip Code**: The provider's zip code.

- **Provider HRR**: The Hospital Referral Region (HRR) where the provider is located.

- **Total Discharges**: The number of discharges billed by the provider for inpatient hospital services.

- **Average Covered Charges**: The provider's average charge for services covered by Medicare for all discharges in the MS-DRG. These will vary from hospital to hospital because of the differences in hospital charge structures.

- **Average Total Payments**: The average total payments to all providers for the MS-DRG including the MSDRG amount, teaching, disproportionate share, capital, and outlier payments for all cases. Also included in the average total payments are co-payment and deductible amounts that the patient is responsible for and any additional payments by third parties for coordination of benefits.

- **Average Medicare Payments**: The average amount that Medicare pays to the provider for Medicare's share of the MS-DRG. Average Medicare payment amounts include the MS-DRG amount, teaching, disproportionate share, capital, and outlier payments for all cases. Medicare payments DO NOT include beneficiary c

# Load file into spark

//StructType objects define the schema of Spark DataFrames. StructType objects contain a list of StructField objects that define the name, type, and nullable flag for each column in a DataFrame.

```scala
package Usecase

import org.apache.spark.sql.SparkSession

import org.apache.spark.sql.types._

object HospitalData_Anlysis {
```

**//StructType objects define the schema of Spark DataFrames. StructType objects contain a list of StructField objects that define the name, type, and nullable flag for each column in a DataFrame.**

```scala
  val CustomSchemaHospital = new StructType(Array(

    StructField("DRGDefinition", StringType,true),

    StructField("ProviderId", LongType,true),

    StructField("ProviderName", StringType,true),

    StructField("ProviderStreetAddress", StringType,true),

    StructField("ProviderCity", StringType,true),

    StructField("ProviderState", StringType,true),

    StructField("ProviderZipCode", LongType,true),

    StructField("HospitalReferralRegionDescription", StringType,true),

    StructField("TotalDischarges", LongType,true),

    StructField("AverageCoveredCharges", DoubleType,true),

    StructField("AverageTotalPayments", DoubleType,true),

    StructField("AverageMedicarePayments", DoubleType,true)))
```

CustomSchemaHospital.printTreeString()

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

  println("hey scala")

  //create a spark session object

  val spark =SparkSession

    .builder()

    .master(master="local")

    .appName(name="Hospital_data_Use_case")

    .config("spark.some.config.option","some-value")

    .getOrCreate()

  println("spark session created")

  spark.sparkContext.setLogLevel("WARN")

  val data
=spark.read.format("csv").option("header",true).schema(CustomSchemaHospital).load("C:/Users/mypc/Desktop/Hospital_data_analysis1.csv").toDF()

  println("Hospital_data_analysis data-->"+data.count())

  data.createOrReplaceTempView("hospital_data")
```

To see the contents inside the DataFrame, type the following:

```scala
data.show()
```

 We will save the data in a table by registering it in a temp table as shown below.

data.createOrReplaceTempView("hospital_data")

## What is the average amount of AverageCoveredCharges per state

**//Objective -1   What is the average amount of AverageCoveredCharges per state???**

```
   val a2= spark.sql("select ProviderState , avg(AverageCoveredCharges) as
Avg_Amount_state from hospital_data group by ProviderState")
```

**//find out the AverageTotalPayments charges per state**

val a3=spark.sql("select ProviderState, avg(AverageTotalPayments) as Avg_Total_Payments from hospital_data group by ProviderState")
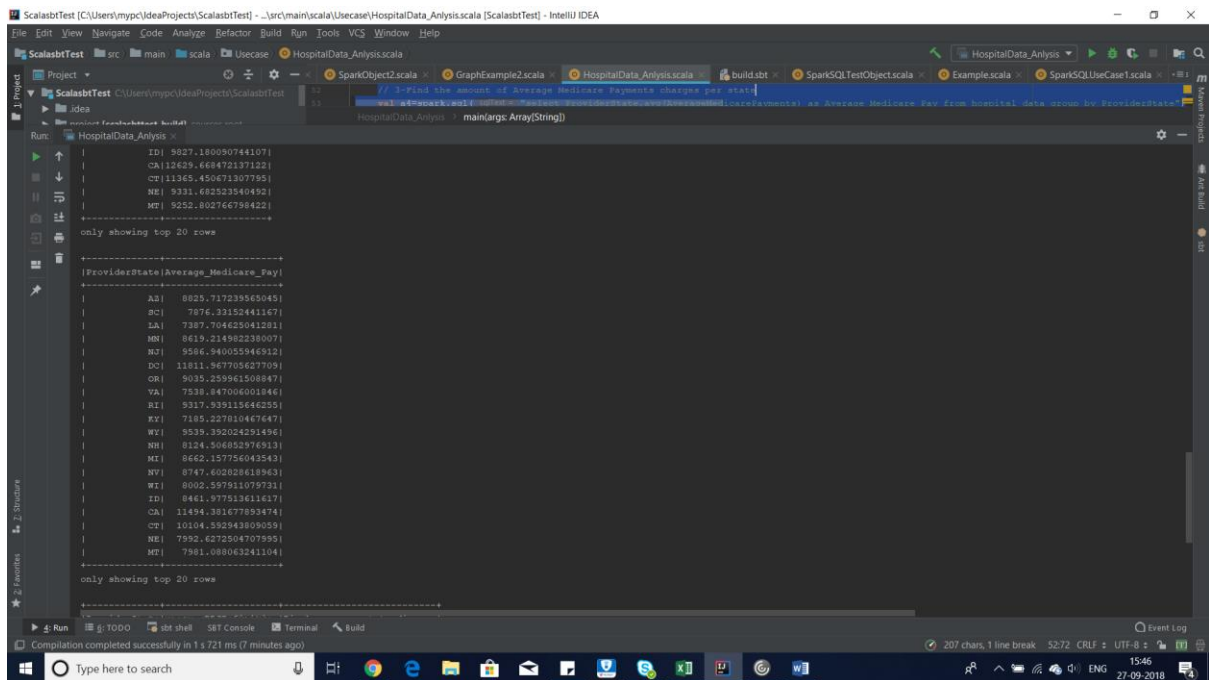
   a3.show()



**// 3-Find the amount of Average Medicare Payments charges per state**

val a4=spark.sql("select ProviderState,avg(AverageMedicarePayments) as Average_Medicare_Pay from hospital_data group by ProviderState")



//4-Find out the total number of Discharges per state and for each disease.

val a5= spark.sql("select ProviderState,DRGDefinition, sum(TotalDischarges)as Discharge_per_state_disease  from hospital_data group by ProviderState,DRGDefinition order by Discharge_per_state_disease desc ")

a5.show()

```
|            NY|  8747.602028618963|
|            WI|  8002.597911075731|
|            ID|  8461.977513611617|
|            CA| 11494.381677893474|
|            CT| 10104.592943809059|
|            NE|  7992.6272504707995|
|            MT|  7981.088063241104|
+--------------+--------------------+
only showing top 20 rows


+--------------+--------------------+-----------------------+
|ProviderState|       DRGDefinition|Discharge_per_state_disease|
+--------------+--------------------+-----------------------+
|            CA|871 - SEPTICEMIA ...|                  34284|
|            TX|470 - MAJOR JOINT...|                  30095|
|            FL|470 - MAJOR JOINT...|                  29985|
|            CA|470 - MAJOR JOINT...|                  29731|
|            TX|871 - SEPTICEMIA ...|                  23144|
|            NY|871 - SEPTICEMIA ...|                  21970|
|            FL|392 - ESOPHAGITIS...|                  21298|
|            IL|470 - MAJOR JOINT...|                  20095|
|            MY|470 - MAJOR JOINT...|                  19371|
|            FL|871 - SEPTICEMIA ...|                  18660|
|            TX|690 - KIDNEY & UR...|                  17384|
|            NY|392 - ESOPHAGITIS...|                  17337|
|            MI|470 - MAJOR JOINT...|                  16847|
|            PA|470 - MAJOR JOINT...|                  16712|
|            FL|292 - HEART FAILU...|                  16639|
|            FL|690 - KIDNEY & UR...|                  16405|
|            OH|470 - MAJOR JOINT...|                  16062|
|            NC|470 - MAJOR JOINT...|                  15820|
|            IL|871 - SEPTICEMIA ...|                  15610|
|            MI|871 - SEPTICEMIA ...|                  15548|
+--------------+--------------------+-----------------------+
only showing top 20 rows


Process finished with exit code 0
```