# Case Study - IV_Hospital_Analysis_in_US

**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**
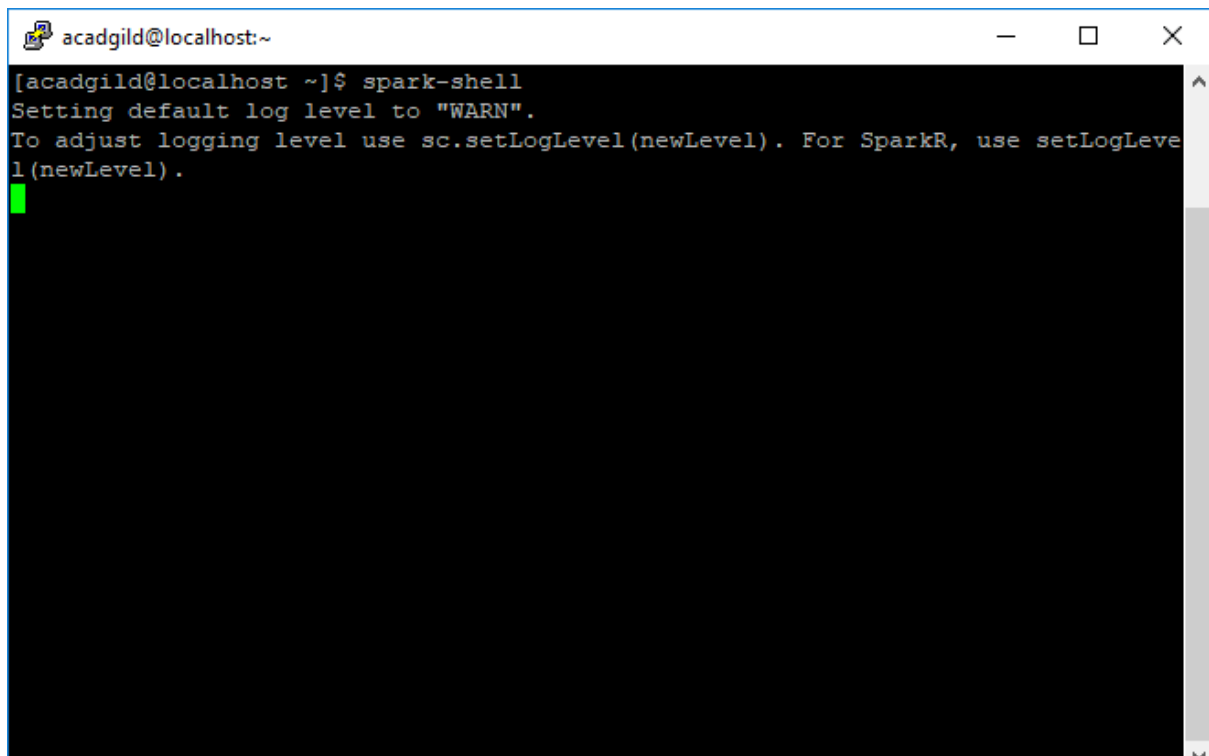
include beneficiary co-payments and deductible amounts nor any additional payments from third parties for coordination of benefits.

**Objective – 1**

**Load file into spark**

**Step 1:**

**Start Spark-shell**



**Step 2:**

**//Load File into Spark**

**val sqlContext = new org.apache.spark.sql.SQLContext(sc)**

**val data = sc.textFile("/user/acadgild/inpatientCharges.csv")**

**//Remove Header**

**val header = data.first()**

**val data1 = data.filter(row => row != header)**

**//Define case class Hospital**

**case class hospital(DRGDefinition:String,ProviderId:String,ProviderName:String,Provide**

rStreetAddress:String,ProviderCity:String,ProviderState:String,ProviderZipCode:String,HospitalReferralRegionDescription:String,TotalDischarges:String,AverageCoveredCharges:String,AverageTotalPayments:String,AverageMedicarePayments:String)

//Convert to DataFrames

val dataDF = data1.map(_.split(",")).map(h
=>hospital(h(0),h(1),h(2),h(3),h(4),h(5),h(6),h(7),h(8),h(9),h(10),h(11))).toDF

```
scala> val sqlContext = new org.apache.spark.sql.SQLContext(sc)
warning: there was one deprecation warning; re-run with -deprecation for details
sqlContext: org.apache.spark.sql.SQLContext = org.apache.spark.sql.SQLContext@2f8739c0

scala> val data = sc.textFile("/user/acadgild/inpatientCharges.csv")
data: org.apache.spark.rdd.RDD[String] = /user/acadgild/inpatientCharges.csv MapPartitionsRDD[360] at textFile at <console>:24

scala> val header = data.first()
header: String = DRGDefinition,ProviderId,ProviderName,ProviderStreetAddress,ProviderCity,ProviderState,ProviderZipCode,HospitalReferralRegionDescription,TotalDischarge
s,AverageCoveredCharges,AverageTotalPayments,AverageMedicarePayments

scala> val data1 = data.filter(row => row != header)
data1: org.apache.spark.rdd.RDD[String] = MapPartitionsRDD[361] at filter at <console>:28

scala> case class hospital(DRGDefinition:String,ProviderId:String,ProviderName:String,ProviderStreetAddress:String,ProviderCity:String,ProviderState:String,ProviderZipC
ode:String,HospitalReferralRegionDescription:String,TotalDischarges:String,AverageCoveredCharges:String,AverageTotalPayments:String,AverageMedicarePayments:String)
defined class hospital

scala> val dataDF = data1.map(_.split(",")).map(h =>hospital(h(0),h(1),h(2),h(3),h(4),h(5),h(6),h(7),h(8),h(9),h(10),h(11))).toDF
dataDF: org.apache.spark.sql.DataFrame = [DRGDefinition: string, ProviderId: string ... 10 more fields]

scala>
```

Step 3:

//Loading Hospital.csv file into temporary table

dataDF.registerTempTable("hospitalTable")

```
scala> dataDF.registerTempTable("hospitalTable")
warning: there was one deprecation warning; re-run with -deprecation for details

scala>
```

**Step 4:**

val result = sqlContext.sql("select * from hospitalTable")

result.show()



**Objective – 2**

1) What is the average amount of AverageCoveredCharges per state

val result1 = sqlContext.sql("select ProviderState,avg(AverageCoveredCharges) as AverageCoveredCharges from hospitalTable GROUP BY ProviderState")

result1.show()

**2) find out the AverageTotalPayments charges per state**

val result2 = sqlContext.sql("select ProviderState,SUM(AverageTotalPayments) AS AverageTotalPayments from hospitalTable GROUP BY ProviderState")

result2.show()

**3) find out the AverageMedicarePayments charges per state.**

val result3 = sqlContext.sql("select
ProviderState,SUM(AverageMedicarePayments) AS
AverageMedicarePayments  from hospitalTable GROUP BY ProviderState")

result3.show()

```
acadgild@localhost:~                                                                                    —    □    ×

scala> val result3 = sqlContext.sql("select ProviderState,SUM(AverageMedicarePayments) AS AverageMedicarePayments  from hospitalTable GROUP BY ProviderState")
result3: org.apache.spark.sql.DataFrame = [ProviderState: string, AverageMedicarePayments: double]

scala> result3.show()
18/04/30 19:28:49 WARN executor.Executor: Managed memory leak detected; size = 17039360 bytes, TID = 12
+-------------------+----------------------+
|      ProviderState|AverageMedicarePayments|
+-------------------+----------------------+
|            TOWANDA|      85933.35999999999|
|          SAN PABLO|      55995.32000000001|
|      PO BOX 1727"|      146123.02000000002|
|         CUMBERLAND|              76058.12|
|            HANCOCK|              16149.57|
|          PRINCETON|      69283.06000000001|
|          WATERTOWN|              65885.64|
|            EDMONDS|      59794.41999999999|
|        MCMINNVILLE|              25478.34|
|               BOAZ|      12612.869999999999|
|             BAXLEY|               4328.57|
|              30002|                  null|
|             140082|                  null|
|             150089|                  null|
|             330024|                  null|
|   750 MORPHY AVENUE|              29557.27|
|2500 ROCKY MOUNTA...|              26560.64|
|         20 YORK ST|              58816.36|
|     1000 MAR-WALT DR|               92185.2|
|     14000 FIVAY ROAD|              65990.82|
+-------------------+----------------------+
only showing top 20 rows


scala>
```
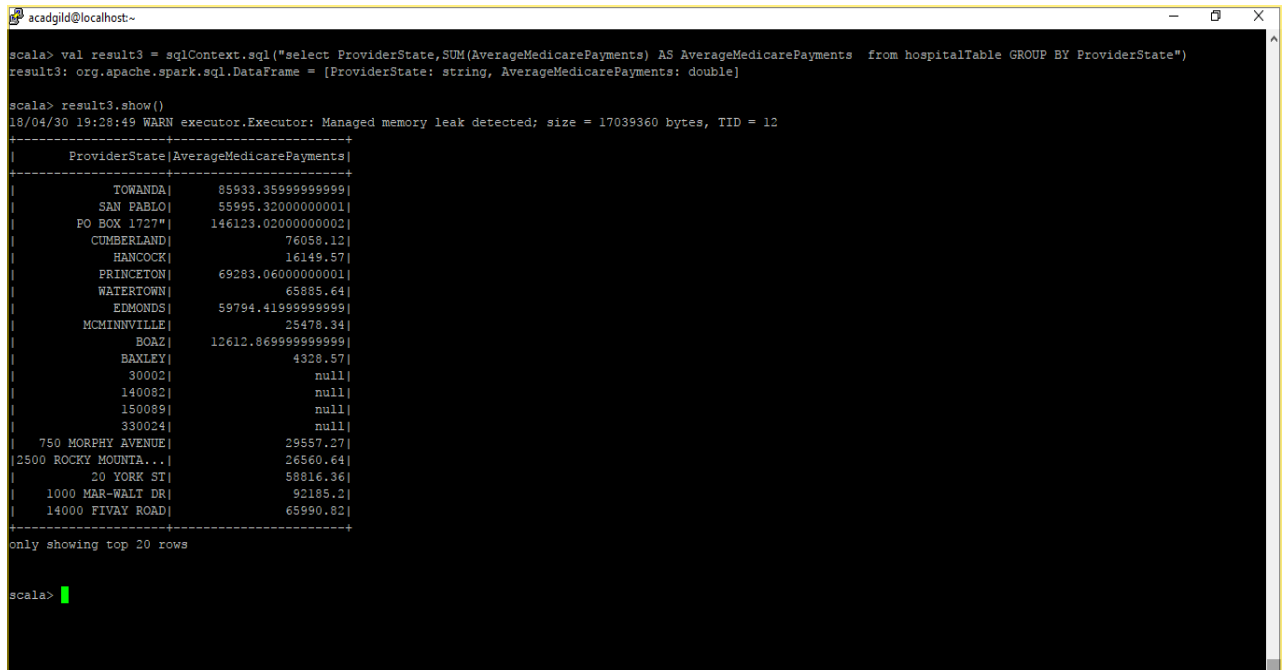
**4) Find out the total number of Discharges per state and for each disease**

val result4 = sqlContext.sql("select
DRGDefinition,ProviderState,COUNT(TotalDischarges) as TotalDischarges
from hospitalTable GROUP BY DRGDefinition,ProviderState")

result4.show()

```
scala> val result4 = sqlContext.sql("select DRGDefinition,ProviderState,COUNT(TotalDischarges) as TotalDischarges from hospitalTable GROUP BY DRGDefinition,ProviderStat
e")
result4: org.apache.spark.sql.DataFrame = [DRGDefinition: string, ProviderState: string ... 1 more field]

scala> result4.show()
18/04/30 19:32:25 WARN executor.Executor: Managed memory leak detected; size = 17039360 bytes, TID = 18
+-------------------+-------------+---------------+
|      DRGDefinition|ProviderState|TotalDischarges|
+-------------------+-------------+---------------+
|057 - DEGENERATIV...|       BANGOR|              1|
|064 - INTRACRANIA...|           AR|             16|
|065 - INTRACRANIA...|  LITTLE ROCK|              1|
|069 - TRANSIENT I...|       BANGOR|              1|
|069 - TRANSIENT I...|  NORTHAMPTON|              1|
|069 - TRANSIENT I...|           OH|             88|
|069 - TRANSIENT I...|      INDIANA|              1|
|189 - PULMONARY E...|           OK|             27|
|192 - CHRONIC OBS...|        DC031|              1|
|193 - SIMPLE PNEU...|           KS|             29|
|193 - SIMPLE PNEU...|    LAFAYETTE|              1|
|193 - SIMPLE PNEU...|    BALTIMORE|              1|
|194 - SIMPLE PNEU...|      MADISON|              1|
|208 - RESPIRATORY...|           ME|              6|
|208 - RESPIRATORY...|           SD|              3|
|238 - MAJOR CARDI...|   GREENSBORO|              1|
|244 - PERMANENT C...|   GREENVILLE|              1|
|254 - OTHER VASCU...|           TX|             66|
|"280 - ACUTE MYOC...|     WHITTIER|              2|
|"280 - ACUTE MYOC...|    LAKE CITY|              2|
+-------------------+-------------+---------------+
only showing top 20 rows


scala>
```

**5) Sort the output in descending order of totalDischarges**

**val result5 = sqlContext.sql("select totalDischarges from hospitalTable SORT BY totalDischarges DESC")**

**result5.show()**

```
scala> val result5 = sqlContext.sql("select totalDischarges from hospitalTable SORT BY totalDischarges DESC")
result5: org.apache.spark.sql.DataFrame = [totalDischarges: string]

scala> result5.show()
+---------------+
|totalDischarges|
+---------------+
|     ZANESVILLE|
|           YUMA|
|     YOUNGSTOWN|
|           YORK|
|        YONKERS|
|        YONKERS|
|         YAKIMA|
|      WYNNEWOOD|
|      WYANDOTTE|
|    WY - Casper|
|    WY - Casper|
|    WY - Casper|
|    WY - Casper|
|    WY - Casper|
|    WY - Casper|
|    WY - Casper|
|    WY - Casper|
| WV - Morgantown|
| WV - Morgantown|
+---------------+
only showing top 20 rows


scala>
```