# Hospital Data

## Code:

**package** spark.basic.cl

**import** org.apache.spark.sql.SparkSession
**import** org.apache.spark.sql.types.IntegerType
**import** org.apache.spark.sql.functions._

**object** Hospital **extends** App{

```
val sparkSession = SparkSession.builder.master("local")
  .appName("spark").getOrCreate()

val sparkcontext = sparkSession.sparkContext

//OBJECTIVE 1

val hp = sparkSession.read.format("csv").option("header","true")
  .load("/home/acadgild/Downloads/inpatientCharges.csv")

hp.show()

import sparkSession.implicits._

//Objective 2

val avgCC = hp.groupBy($"ProviderState").agg(avg("AverageCoveredCharges")).show()

val sumATP = hp.groupBy($"ProviderState").agg(sum("AverageTotalPayments")).show()

val sumAMP =
hp.groupBy($"ProviderState").agg(sum("AverageMedicarePayments")).show()

//OBJECTIVE 3

val PSD = hp.groupBy($"ProviderState",$"DRGDefinition").agg(sum("TotalDischarges"))
PSD.show()

val DEOR = PSD.sort(desc("sum(TotalDischarges)")).show()

}
```

### OBJECTIVE 1

**Load file into spark**

```
val hp = sparkSession.read.format("csv").option("header","true")
  .load("/home/acadgild/Downloads/inpatientCharges.csv")

hp.show()
```

## OBJECTIVE 2

➤ **What is the average amount of AverageCoveredCharges per state**

```scala
val avgCC = hp.groupBy($"ProviderState").agg(avg("AverageCoveredCharges")).show()
```

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

**val** *sumATP* = *hp*.groupBy(**$"ProviderState"**).agg(*sum*(**"AverageTotalPayments"**)).show()



➤ **find out the AverageMedicarePayments charges per state.**

**val** *sumAMP* =
*hp*.groupBy(**$"ProviderState"**).agg(*sum*(**"AverageMedicarePayments"**)).show()

## OBJECTIVE 3

➤ **Find out the total number of Discharges per state and for each disease**

**val** *PSD* = *hp*.groupBy(**$"ProviderState",$"DRGDefinition"**).agg(*sum*(**"TotalDischarges"**))
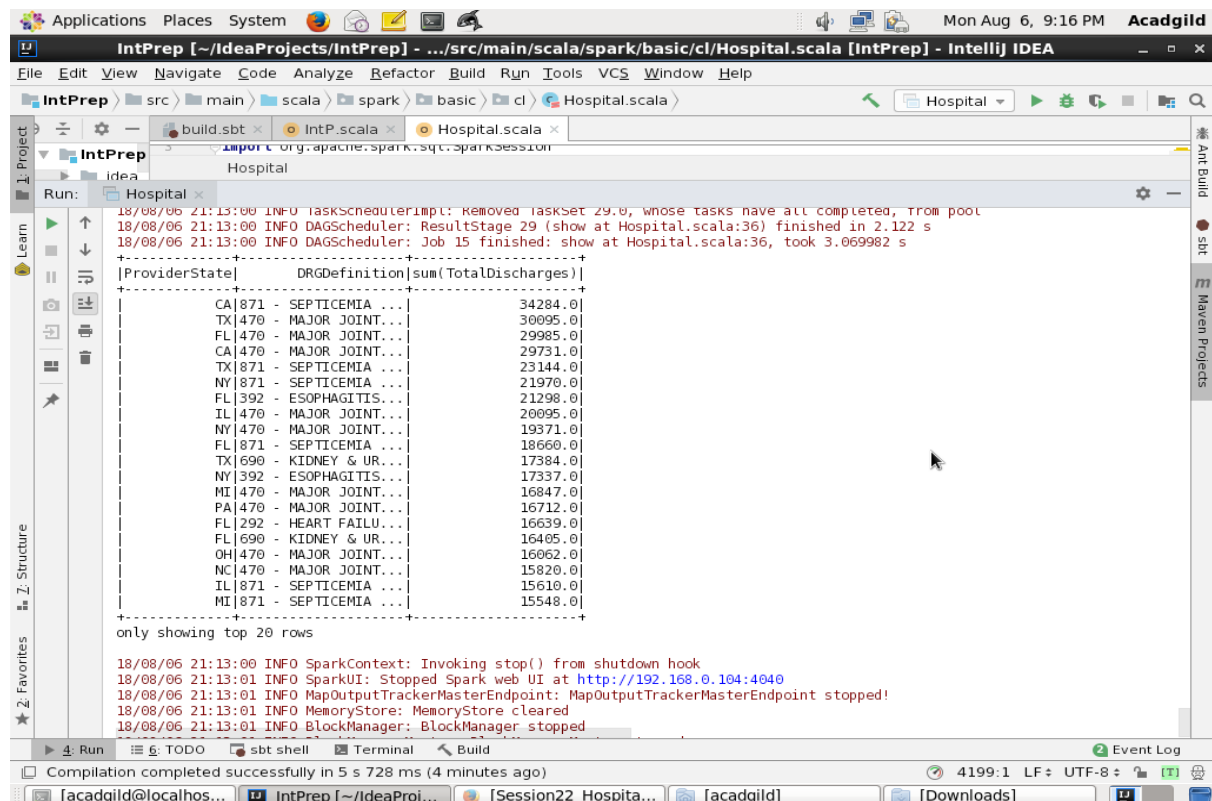  *PSD*.show()



➤ **Sort the output in descending order of totalDischarges**

**val** *DEOR* = *PSD*.sort(*desc*(**"sum(TotalDischarges)"**)).show()