

Case Study 4

```
import org.apache.spark.sql.Session
import org.apache.spark.sql.functions._
object CaseStudy4 {

  case class Hospital(DRGDefinition: String, ProviderId: Int, ProviderName: String, street:
String, city: String, state: String, zip: Int, HReferra: String, TotalDischarges: Int,
avg_cover_charge: Double, avg_tot_pay: Double, avg_medicare: Double)

  def main(args: Array[String]): Unit = {
    println("hey scala")
    val spark = Session
      .builder()
      .master("local")
      .appName("Case study 4")
      .config("spark.some.config.option", "some-value")
      .getOrCreate()
    println("Spark Session Object created")

    val hosp = spark.sqlContext.read.csv("E:\\Avani\\Acadgild\\Case Study
4\\inpatientCharges.csv")
    val hosDF= hosp.toDF()
    hosDF.show(truncate = false)

    hosDF.registerTempTable("medical")
  }
}
```

Objective 1

- Load file into spark

```
val hosp = spark.sqlContext.read.csv("E:\\Avani\\Acadgild\\Case Study
4\\inpatientCharges.csv")
val hosDF= hosp.toDF()
hosDF.show(truncate = false)

hosDF.registerTempTable("medical")
```

Objective2

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

```
val avg_cover=spark.sql("SELECT (_c5),AVG(_c9) FROM MEDICAL GROUP BY (_c5)")
avg_cover.show()
println("Average covered charges")
```

- Find out the **AverageTotalPayments** charges per state

```
val avg_payment=spark.sql("SELECT (_c5),sum(_c10) FROM MEDICAL GROUP BY (_c5) order
by (_c5)")

avg_payment.show()

println("Average Total payment charges")
```

- Find out the **AverageMedicarePayments** charges per state

```
val avg_med=spark.sql("SELECT (_c5),sum(_c11) FROM MEDICAL GROUP BY (_c5) order by
(_c5)")

avg_med.show()

println("Average Medicare payment charges")
```

Objective3

- Find out the total number of **Discharges** per state and for each disease
- Sort the output in descending order of **totalDischarges**

```
val tot_dis=spark.sql("SELECT (_c5), (_c0), sum(_c8) total_discharge FROM MEDICAL GROUP
BY (_c5), (_c0)")

tot_dis.show(truncate = false)    //to display the complete content of a column

println("Total discharges")

val order= tot_dis.orderBy(desc("total_discharge"))    //sorting the above data in
descending order

order.show(truncate = false)

println("sorted data")
```

Screenshots

```
18/06/05 22:31:31 INFO BlockManagerInfo: Removed broadcast_0_piece0 on 192.168.5.112:61238 in memory (size: 14.3 KB, free: 899.7 MB)
18/06/05 22:31:31 INFO Executor: Finished task 0.0 in stage 1.0 (TID 1). 4467 bytes result sent to driver
18/06/05 22:31:31 INFO TaskSchedulerImpl: Removed TaskSet 1.0, whose tasks have all completed, from pool
18/06/05 22:31:31 INFO DAGScheduler: ResultStage 1 (show at CaseStudy4.scala:21) finished in 0.257 s
18/06/05 22:31:31 INFO DAGScheduler: Job 1 finished: show at CaseStudy4.scala:21, took 0.332478 s
18/06/05 22:31:31 INFO CodeGenerator: Code generated in 56.028994 ms
```

l_c0	l_c1	l_c2	l_c3	l_c4	l_c5	l_c6	l_c7
DRGDefinition	ProviderId	ProviderName	ProviderStreetAddress	ProviderCity	ProviderState	ProviderZipCode	HospitalReferr
039 - EXTRACRANIAL PROCEDURES W/O CC/MCC10001		SOUTHEAST ALABAMA MEDICAL CENTER	1108 ROSS CLARK CIRCLE	DOTHAN	AL	36301	AL - Dothan
039 - EXTRACRANIAL PROCEDURES W/O CC/MCC10005		MARSHALL MEDICAL CENTER SOUTH	2505 U S HIGHWAY 431 NORTH	BOAZ	AL	35957	AL - Birmingham
039 - EXTRACRANIAL PROCEDURES W/O CC/MCC10006		ELIZA COFFEE MEMORIAL HOSPITAL	205 MARENGO STREET	FLORENCE	AL	35631	AL - Birmingham
039 - EXTRACRANIAL PROCEDURES W/O CC/MCC10011		ST VINCENT'S EAST	150 MEDICAL PARK EAST DRIVE	BIRMINGHAM	AL	35235	AL - Birmingham
039 - EXTRACRANIAL PROCEDURES W/O CC/MCC10016		SHELBY BAPTIST MEDICAL CENTER	1000 FIRST STREET NORTH	ALABASTER	AL	35007	AL - Birmingham
039 - EXTRACRANIAL PROCEDURES W/O CC/MCC10023		BAPTIST MEDICAL CENTER SOUTH	2105 EAST SOUTH BOULEVARD	MONTGOMERY	AL	36116	AL - Montgomer
039 - EXTRACRANIAL PROCEDURES W/O CC/MCC10029		EAST ALABAMA MEDICAL CENTER AND SNF	2000 PEEPERELL PARKWAY	OPELIKA	AL	36801	AL - Birmingham
039 - EXTRACRANIAL PROCEDURES W/O CC/MCC10033		UNIVERSITY OF ALABAMA HOSPITAL	619 SOUTH 16TH STREET	BIRMINGHAM	AL	35333	AL - Birmingham
039 - EXTRACRANIAL PROCEDURES W/O CC/MCC10038		HUNTSVILLE HOSPITAL	101 SIVLEY RD	HUNTSVILLE	AL	35801	AL - Huntsville
039 - EXTRACRANIAL PROCEDURES W/O CC/MCC10040		GADSDEN REGIONAL MEDICAL CENTER	1007 GOODYEAR AVENUE	GADSDEN	AL	35903	AL - Huntsville
039 - EXTRACRANIAL PROCEDURES W/O CC/MCC10046		RIVERSVIEW REGIONAL MEDICAL CENTER	1600 SOUTH THIRD STREET	GADSDEN	AL	35901	AL - Birmingham
039 - EXTRACRANIAL PROCEDURES W/O CC/MCC10055		FLOWERS HOSPITAL	14370 WEST MAIN STREET	DOTHAN	AL	36305	AL - Dothan
039 - EXTRACRANIAL PROCEDURES W/O CC/MCC10056		ST VINCENT'S BIRMINGHAM	810 ST VINCENT'S DRIVE	BIRMINGHAM	AL	35205	AL - Birmingham
039 - EXTRACRANIAL PROCEDURES W/O CC/MCC10078		NORTHEAST ALABAMA REGIONAL MED CENTER	400 EAST 10TH STREET	ANNISTON	AL	36207	AL - Birmingham
039 - EXTRACRANIAL PROCEDURES W/O CC/MCC10083		SOUTH BALDWIN REGIONAL MEDICAL CENTER	1613 NORTH MCKENZIE STREET	FOLEY	AL	36535	AL - Mobile
039 - EXTRACRANIAL PROCEDURES W/O CC/MCC10085		DECATUR GENERAL HOSPITAL	1201 7TH STREET SE	DECATUR	AL	35609	AL - Huntsville
039 - EXTRACRANIAL PROCEDURES W/O CC/MCC10090		PROVIDENCE HOSPITAL	16801 AIRPORT BOULEVARD	MOBILE	AL	36608	AL - Mobile
039 - EXTRACRANIAL PROCEDURES W/O CC/MCC10092		D C H REGIONAL MEDICAL CENTER	809 UNIVERSITY BOULEVARD EAST	TUSCALOOSA	AL	35401	AL - Tuscaloosa
039 - EXTRACRANIAL PROCEDURES W/O CC/MCC10100		THOMAS HOSPITAL	1750 MORPHY AVENUE	FAIRHOPE	AL	36532	AL - Mobile

only showing top 20 rows

```
18/06/05 22:31:31 INFO SparkContext: Invoking stop() from shutdown hook
```

1.a

```
18/06/05 22:31:31 INFO BlockManagerInfo: Removed broadcast_0_piece0 on 192.168.5.112:61238 in memory (size: 14.3 KB, free: 899.7 MB)
18/06/05 22:31:31 INFO Executor: Finished task 0.0 in stage 1.0 (TID 1). 4467 bytes result sent to driver
18/06/05 22:31:31 INFO TaskSchedulerImpl: Removed TaskSet 1.0, whose tasks have all completed, from pool
18/06/05 22:31:31 INFO DAGScheduler: ResultStage 1 (show at CaseStudy4.scala:21) finished in 0.257 s
18/06/05 22:31:31 INFO DAGScheduler: Job 1 finished: show at CaseStudy4.scala:21, took 0.332478 s
18/06/05 22:31:31 INFO CodeGenerator: Code generated in 56.028994 ms
```

l_c4	l_c5	l_c6	l_c7	l_c8	l_c9	l_c10	l_c11	
Address	ProviderCity	ProviderState	ProviderZipCode	HospitalReferralRegionDescription	TotalDischarges	AverageCoveredCharges	AverageTotalPayments	AverageMedicarePayments
K CIRCLE	DOTHAN	AL	36301	AL - Dothan	191	132963.07	15777.24	14763.73
AY 431 NORTH	BOAZ	AL	35957	AL - Birmingham	114	115131.85	15787.57	14976.71
REET	FLORENCE	AL	35631	AL - Birmingham	124	137560.37	15434.95	14453.79
K EAST DRIVE	BIRMINGHAM	AL	35235	AL - Birmingham	125	113999.28	15117.56	14129.16
EET NORTH	ALABASTER	AL	35007	AL - Birmingham	118	131633.27	15659.38	14851.44
H BOULEVARD	MONTGOMERY	AL	36116	AL - Montgomery	167	116920.79	16652.8	15374.14
PARKWAY	OPELIKA	AL	36801	AL - Birmingham	151	111977.13	15834.74	14761.41
STREET	BIRMINGHAM	AL	35233	AL - Birmingham	132	135841.09	18031.12	15859.5
	HUNTSVILLE	AL	35801	AL - Huntsville	135	128523.39	16113.38	15228.4
AVENUE	GADSDEN	AL	35903	AL - Birmingham	134	175233.38	15541.05	14386.94
D STREET	GADSDEN	AL	35901	AL - Birmingham	114	167327.92	15461.57	14493.57
STREET	DOTHAN	AL	36305	AL - Dothan	145	139607.28	15356.28	14408.2
'S DRIVE	BIRMINGHAM	AL	35205	AL - Birmingham	143	122862.23	15374.65	14186.02
STREET	ANNISTON	AL	36207	AL - Birmingham	121	131110.85	15366.23	14376.23
ENZIE STREET	FOLEY	AL	36535	AL - Mobile	115	125411.33	15282.93	14383.73
T SE	DECATUR	AL	35609	AL - Huntsville	127	19234.51	15676.55	14509.11
BOULEVARD	MOBILE	AL	36608	AL - Mobile	127	115895.85	15930.11	13972.85
BOULEVARD EAST	TUSCALOOSA	AL	35401	AL - Tuscaloosa	131	119721.16	16192.54	15179.38
NOE	FAIRHOPE	AL	36532	AL - Mobile	118	10710.88	14968	13896.88

only showing top 20 rows

```
18/06/05 22:31:31 INFO SparkContext: Invoking stop() from shutdown hook
```

1.b(contd..)

```
18/06/06 19:53:30 INFO DAGScheduler: ResultStage 9 (show at CaseStudy4.scala:31) finished in 1.432 s
18/06/06 19:53:30 INFO DAGScheduler: Job 5 finished: show at CaseStudy4.scala:31, took 1.484162 s
18/06/06 19:53:30 INFO CodeGenerator: Code generated in 18.851803 ms

+-----+
|_c5|avg(CAST(_c9 AS DOUBLE))|
+-----+
|AZ|41200.063019992595|
|SC|35862.48456269756|
|LA|33085.372791542646|
|MI|27894.36182060388|
|NY|66125.6862743729|
|DC|40116.66365800864|
|OR|27390.111870669723|
|VA|29222.000487072903|
|RI|29942.701122448976|
|KY|24523.80716940223|
|WY|28700.59862348178|
|NH|27059.020801944105|
|ME|24124.247209817277|
|NV|61047.11541597937|
|WJ|26149.828331686607|
|ID|25565.547041742288|
|CA|67508.616535517|
|CT|31318.4101143709|
|NE|31736.427824858758|
|MT|22670.015237154144|
+-----+

only showing top 20 rows

Average covered charges
18/06/06 19:53:30 INFO SparkContext: Invoking stop() from shutdown hook
18/06/06 19:53:30 INFO BlockManagerInfo: Removed broadcast 9 stored on 192.168.0.115:63832 (18.6 MB, free: 888.6 MB)

Compilation completed successfully with 1 warning in 2 s 959 ms (moments ago)
```

2.a

```
18/06/06 19:54:45 INFO DAGScheduler: Job 2 finished: show at CaseStudy4.scala:35, took 4.821457 s

+-----+
|_c5|sum(CAST(_c10 AS DOUBLE))|
+-----+
|AK|3366222.489999999|
|AL|2.751052385999973E7|
|AR|1.6575787280000016E7|
|AZ|2.8950559930000026E7|
|CA|1.6499398891999936E8|
|CO|1.796007568999996E7|
|CT|2.285592129999975E7|
|DE|6005089.589999995|
|FL|4081868.5299999984|
|GA|9.846807830999967E7|
|HI|4.434384416999986E7|
|IL|5646876.870000002|
|IN|1.441399933999998E7|
|ID|5414776.230000002|
|IA|7.74344572499997E7|
|IH|3.730091159000008E7|
|KS|1.385007038000008E7|
|KY|2.673156338000008E7|
|LA|2.614923161999968E7|
|MA|3.949568905999987E7|
+-----+

only showing top 20 rows

Average Total payment charges
18/06/06 19:54:45 INFO SparkContext: Invoking stop() from shutdown hook
18/06/06 19:54:45 INFO BlockManagerInfo: Removed broadcast 9 stored on 192.168.0.115:63832 (18.6 MB, free: 888.6 MB)

Compilation completed successfully with 1 warning in 3 s 33 ms (a minute ago)
```

2.b

```
10/06/06 19:55:53 INFO DAGScheduler: ResultStage 3 (show at CaseStudy4.scala:35) finished in 3.032 s
10/06/06 19:55:53 INFO DAGScheduler: Job 2 finished: show at CaseStudy4.scala:39, took 4.760928 s
+-----+
|_c5|sum(CAST(_c11 AS DOUBLE))|
+-----+
| AK | 2993521.9399999998 |
| AL | 2.332945587999997E7 |
| AR | 1.4303062960000006E7 |
| AZ | 2.5162119849999946E7 |
| CA | 1.5016260224000034E8 |
| CO | 1.5405260329999983E7 |
| CT | 2.032033641000002E7 |
| DC | 5457129.080000001 |
| DE | 3530111.2699999998 |
| FL | 8.553072483999966E7 |
| GA | 3.809254514000004E7 |
| HI | 4847623.969999999 |
| IA | 1.239484006999998E7 |
| ID | 4662549.610000001 |
| IL | 6.631920065999971E7 |
| IN | 3.185754233999903E7 |
| KS | 1.1833965499999976E7 |
| KY | 2.320110060000003E7 |
| LA | 2.2362581899999958E7 |
| MA | 3.550668563999989E7 |
+-----+
only showing top 20 rows

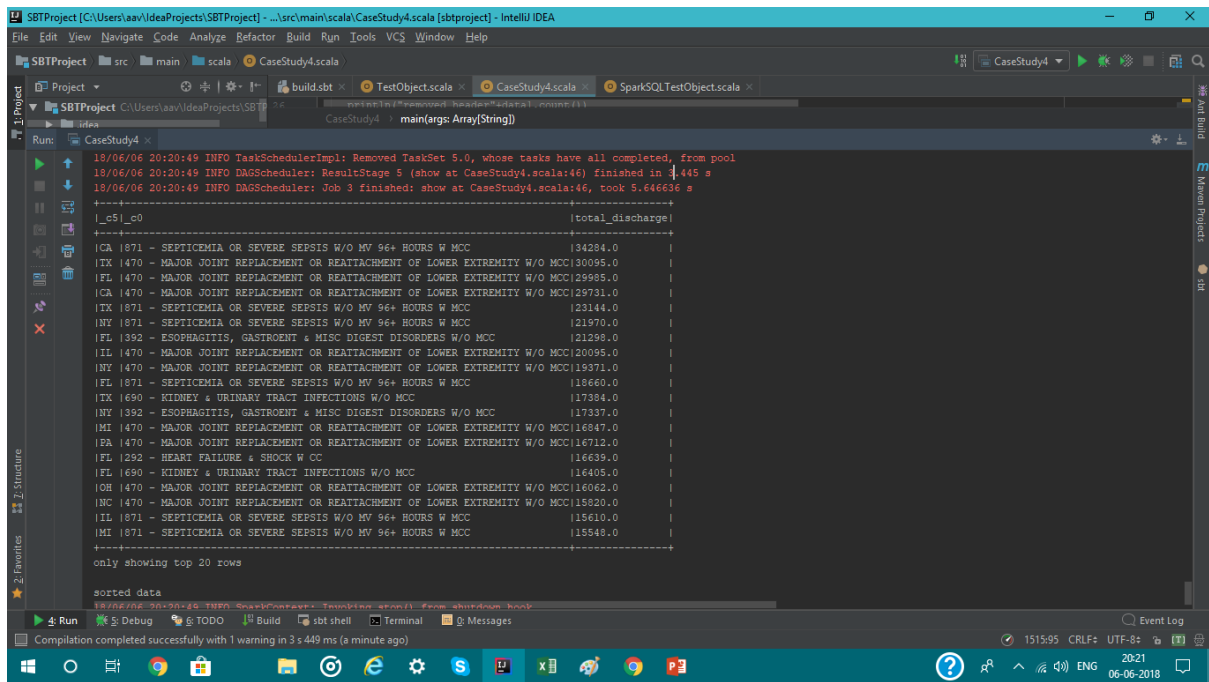
Average Medicare payment charges
10/06/06 19:55:53 INFO CodeGenerator: Code generated in 27.7654 ms
```

2.c

```
10/06/06 20:20:43 INFO CodeGenerator: Code generated in 14.948363 ms
+-----+
|_c5|_c0|total_discharge|
+-----+
| KY | 1065 - INTRACRANIAL HEMORRHAGE OR CEREBRAL INFARCTION W CC | 1937.0 |
| NY | 1101 - SEIZURES W/O MCC | 14503.0 |
| IN | 1149 - DYSEQUILIBRIUM | 1700.0 |
| IA | 1178 - RESPIRATORY INFECTIONS & INFLAMMATIONS W CC | 1540.0 |
| WI | 1202 - BRONCHITIS & ASTHMA W CC/MCC | 1338.0 |
| MO | 1208 - RESPIRATORY SYSTEM DIAGNOSIS W VENTILATOR SUPPORT <96 HOURS | 1840.0 |
| WI | 1251 - PERC CARDIOVASC PROC W/O CORONARY ARTERY STENT W/O MCC | 1417.0 |
| AR | 1261 - ACUTE MYOCARDIAL INFARCTION, DISCHARGED ALIVE W CC | 1413.0 |
| AZ | 1282 - HEART FAILURE & SHOCK W CC | 12643.0 |
| NY | 1282 - HEART FAILURE & SHOCK W CC | 13289.0 |
| NV | 1293 - HEART FAILURE & SHOCK W/O CC/MCC | 519.0 |
| SD | 1303 - ATHEROSCLEROSIS W/O MCC | 53.0 |
| TN | 1305 - HYPERTENSION W/O MCC | 1730.0 |
| ME | 1308 - CARDIAC ARRHYTHMIA & CONDUCTION DISORDERS W MCC | 312.0 |
| NV | 1372 - MAJOR GASTROINTESTINAL DISORDERS & PERITONEAL INFECTIONS W CC | 1126.0 |
| WA | 1392 - ESOPHAGITIS, GASTROENT & MISC DIGEST DISORDERS W/O MCC | 3148.0 |
| WI | 1439 - DISORDERS OF PANCREAS EXCEPT MALIGNANCY W CC | 1215.0 |
| MN | 1536 - FRACTURES OF HIP & PELVIS W/O MCC | 1332.0 |
| DC | 1563 - FX, SPN, STRN & DISL EXCEPT FEMUR, HIP, PELVIS & THIGH W/O MCC | 143.0 |
| CO | 1602 - CELLULITIS W MCC | 186.0 |
+-----+
only showing top 20 rows

Total discharges
10/06/06 20:20:43 INFO FileSourceStrategy: Pruning directories with:
10/06/06 20:20:43 INFO FileSourceStrategy: Post-Scan Filters:
10/06/06 20:20:43 INFO FileSourceStrategy: Format Data Schema: structure: c0: string, c5: string, c6: string, 1 more fields
```

3.a



3.b

Complete code

```
import org.apache.spark.sql.Session
import org.apache.spark.sql.functions._
object CaseStudy4 {
```

```
  case class Hospital(DRGDefinition: String, ProviderId: Int, ProviderName: String, street:
String, city: String, state: String, zip: Int, HReferra: String, TotalDischarges: Int,
avg_cover_charge: Double, avg_tot_pay: Double, avg_medicare: Double)
```

```
  def main(args: Array[String]): Unit = {
    println("hey scala")
    val spark = SparkSession
      .builder()
      .master("local")
      .appName("Case study 4")
      .config("spark.some.config.option", "some-value")
      .getOrCreate()
    println("Spark Session Object created")
```

```
    val hosp = spark.sqlContext.read.csv("E:\\Avani\\Acadgild\\Case Study
```

```

4\\inpatientCharges.csv")
val hosDF= hosp.toDF()
  hosDF.show(truncate = false)
val header=hosp.first()
val data1 = hosp.filter(row => row != header) //removing of header
data1.count()
println("Hospital Data->>" + hosp.show(truncate = false))
println("removed header" + data1.count())

hosDF.registerTempTable("medical")

val avg_cover=spark.sql("SELECT (_c5),AVG(_c9) FROM MEDICAL GROUP BY (_c5)")
avg_cover.show()
println("Average covered charges")

val avg_payment=spark.sql("SELECT (_c5),sum(_c10) FROM MEDICAL GROUP BY (_c5)
order by (_c5)")
avg_payment.show()
println("Average Total payment charges")

val avg_med=spark.sql("SELECT (_c5),sum(_c11) FROM MEDICAL GROUP BY (_c5) order by
(_c5)")
avg_med.show()
println("Average Medicare payment charges")

val tot_dis=spark.sql("SELECT (_c5), (_c0), sum(_c8) total_discharge FROM MEDICAL
GROUP BY (_c5), (_c0)")
tot_dis.show(truncate = false)
println("Total discharges")
val order= tot_dis.orderBy(desc("total_discharge"))
  order.show(truncate = false)
  println("sorted data")

}
}

```